

Disentangled multi-subject and social behavioral representations through a constrained subspace variational autoencoder (CS-VAE)

Daiyao Yi¹ Simon Musall² Anne Churchland³ Nancy Padilla-Coreano⁴
Shreya Saxena¹

¹Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL, USA.

{yidaiyao, shreya.saxena}@ufl.edu

²Department of Neurophysiology, Institute of Biology 2, RWTH Aachen University, Aachen, Germany

³Department of Neurobiology, University of California, Los Angeles, Los Angeles, CA, USA

⁴Department of Neuroscience, University of Florida, Gainesville, FL, USA

Abstract

Effectively modeling and quantifying behavior is essential for our understanding of the brain. Modeling behavior in naturalistic settings in social and multi-subject tasks remains a significant challenge. Modeling the behavior of different subjects performing the same task requires partitioning the behavioral data into features that are common across subjects, and others that are distinct to each subject. Modeling social interactions between multiple individuals in a freely-moving setting requires disentangling effects due to the individual as compared to social investigations. To achieve flexible disentanglement of behavior into interpretable latent variables with individual and across-subject or social components, we build on a semi-supervised approach to partition the behavioral subspace, and propose a novel regularization based on the Cauchy-Schwarz divergence to the model. Our model, known as the constrained subspace variational autoencoder (CS-VAE), successfully models distinct features of the behavioral videos across subjects, as well as continuously varying differences in social behavior. Our approach vastly facilitates the analysis of the resulting latent variables in downstream tasks such as uncovering disentangled behavioral motifs, the efficient decoding of a novel subject's behavior, and provides an understanding of how similarly different animals perform innate behaviors.

1 Introduction

Effective study of the relationship between neural signals and ensuing behavior relies on our ability to measure and adequately quantify behavior. Historically, behavior has been quantified by a very small number of markers as the subject performs the task, for example, force sensors on levers. However, advancement in hardware and storage capabilities, as well as computational methods applied to video data, has allowed us to increase the quality and capability of behavioral recordings to videos of the entire subject that can be processed and analyzed quickly. It is now widely recognized that understanding the relationship between complex neural activity and high-dimensional behavior is a major step in understanding the brain that has been undervalued in the past [1, 2]. However, the analysis of high-dimensional behavioral video data across subjects is still a nascent field, due to the lack of adequate tools to efficiently disentangle behavioral features related to different subjects. Moreover, as recording modalities become light-weight and portable, neural and behavioral recordings can be performed in more naturalistic settings, which are difficult for behavioral analysis tools to disentangle due to changing scenes.

Although pose estimation tools that track various body parts in a behavioral video are very popular, they fail to capture smaller movements and rely on the labeler to judge which parts of the scene are important to track [3, 4, 5, 6, 7]. Unsupervised techniques have gained traction to circumvent these problems. These

17 include directly applying dimensionality reduction methods such as Principal Component Analysis (PCA) and
18 Variational Autoencoders (VAEs) to video data [2, 8, 9]. However, understanding or segmentation of the latent
19 variables is difficult for any downstream tasks such as motif generation. To combine the best of both worlds,
20 semi-supervised VAEs have been used for the joint estimation of tracked body parts and unsupervised latents
21 that can effectively describe the entire image [2]. These have not been applied to across-subject data, with the
22 exception of [10], where the authors directly use a frame of each subject’s video as a context frame to define
23 individual differences; however, this method only works with a *discrete* set of *labeled* sessions or subjects. These
24 methods fail when applied without labeled subject data, or more importantly, when analyzing freely-behaving
25 social behavior, due to continuously shifting image distributions that confound the latent space.

26 With increasing capabilities to effectively record more naturalistic data in neuroscience, there is a growing
27 demand for behavioral analysis methods that are tailored to these settings. In this work, we model a
28 continuously varying distribution of images, such as in freely moving and multi-subject behavior, by using a
29 novel loss term called the Cauchy-Schwarz Divergence (CSD) [11, 12]. By applying the CSD loss term, a
30 subset of the latents can be automatically projected on a pre-defined and flexible distribution, thus leading
31 to an unbiased approach towards latent separation. Here, the CSD is an effective variational regularizer
32 that separates the latents corresponding to images with different appearances, thus successfully capturing
33 ‘background’ information of an individual. This background information can be the difference in lighting
34 during the experiment, the difference in appearance across mice in a multi-subject dataset, or the presence of
35 another subject in the same field of view as in a social interaction dataset.

36 To further demonstrate the utility of our approach, we show that we can recover behavioral motifs from
37 the resulting latents in a seamless manner. We recover (a) the same motifs across different animals performing
38 the same task, and (b) motifs pertaining to social interactions in a freely moving task with two animals.
39 Furthermore, we show the neural decoding of multiple animals in a unified model, with benefits towards the
40 efficient decoding of the behavior of a novel subject. Finally, we compare the commonalities in neural activity
41 across different trials in the same subject to those across subjects for different types of behavior motifs, e.g.
42 task-related and spontaneous.

43 **Related Works** Pose estimation tools such as DeepLabCut (DLC) and LEAP have been broadly applied
44 to neuroscience experiments to track the body parts of animals performing different tasks, including in the
45 social setting [3, 4, 5, 6, 7]. These are typically supervised techniques that require extensive manual labeling.
46 Although these methods can be sample-efficient due to the use of transfer learning methods, they still depend
47 inherently on the quality of the manual labels, which can differ across labelers. Moreover, these methods may
48 be missing key information in these behavioral videos that are not captured by tracking the body parts, for
49 example, movements of the face, the whiskers, and smaller muscles that comprise a subject’s movements.

50 Emerging unsupervised methods have demonstrated significant potential in directly modeling behavioral
51 videos. A pioneer in this endeavor was MoSeq, a behavioral video analysis tool that encodes high dimensional
52 behavior by directly applying PCA to the data [13, 9]. Behavenet is similar to MoSeq, but uses autoencoders
53 to more effectively reduce the dimensionality of the representation [8]. However, the corresponding latent
54 variables in these models are typically not interpretable. To add interpretability, the Partitioned Subspace
55 VAE (PS-VAE) [2] formulates a semi-supervised approach that uses the labels generated using pose estimation
56 methods such as DLC in order to partition the latent representation into both supervised and unsupervised
57 subspaces. The ‘supervised’ latent subspace captures the parts that are labeled by pose estimation software,
58 while the ‘unsupervised’ latent subspace encodes the parts of the image that have not been accounted for
59 by the supervised space. While PS-VAE is very effective for a single subject, it does not address latent
60 disentanglement in the ‘unsupervised’ latent space, and is not able to model multi-subject or social behavioral
61 data.

62 Modeling multiple sessions has recently been examined in two approaches: MSPS-VAE and DBE [2, 10].
63 Both of these are confined to modeling head-fixed animals with a pre-specified number of sessions or subjects.
64 In MSPS-VAE, an extension to PS-VAE, a latent subspace is introduced in the model that encodes the static
65 differences across sessions. In DBE, a context frame from each session or subject is used as a static input
66 to generate the behavioral embeddings. Two notable requirements of applying both these methods is the
67 presence of a discrete number of labeled sessions or subjects in the dataset. Therefore, these are not well
68 suited for naturalistic settings where the session / subject identity might not be known a priori, or the scene
69 might be continuously varying, for example, in the case of subjects roaming in an open-field.

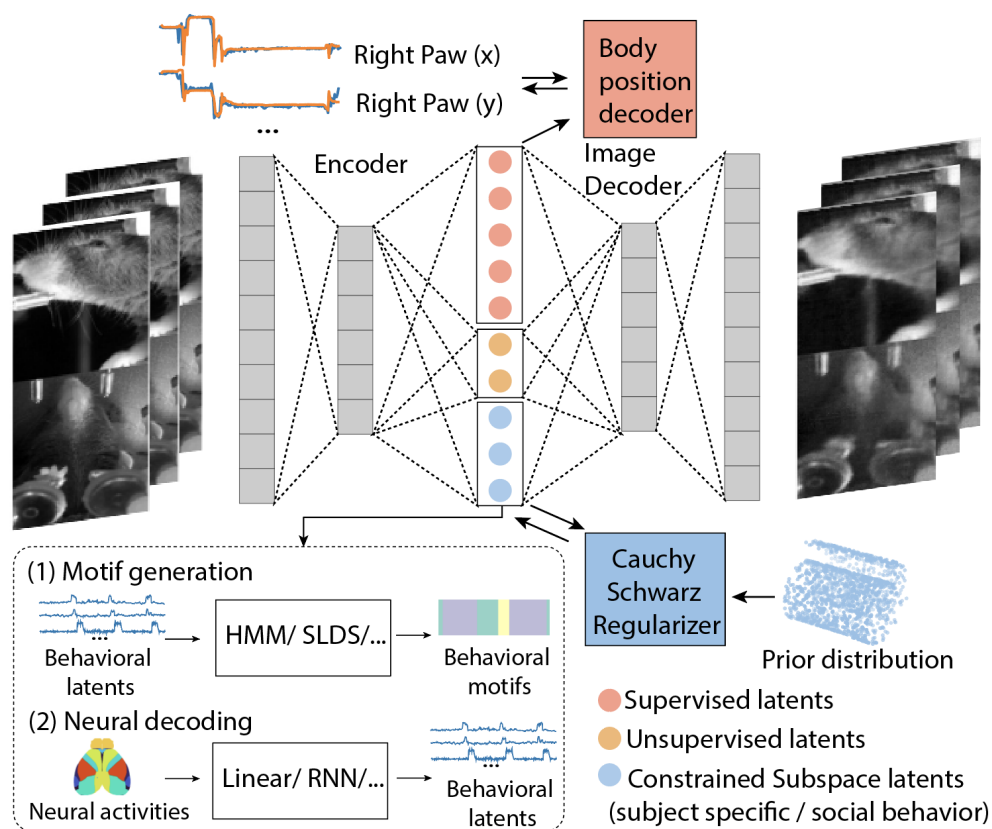


Figure 1: Overview of the Constrained Subspace Variational Autoencoder (CS-VAE). The latent space is divided in three parts: (1) the supervised latents decode the labeled body positions, (2) the unsupervised latents model the individual’s behavior that is not explained by the supervised latents, and (3) the constrained subspace latents model the continuously varying features of the image, e.g., relating to multi-subject or social behavior. After training the network, the generated latents can be applied to several downstream tasks. Here we show two example tasks: (1) Motif generation: we apply state space models such as hidden Markov models (HMM) and switched linear dynamical systems (SLDS), with the behavioral latent variables as the observations; (2) Neural decoding: with neural recordings such as widefield calcium imaging, corresponding behaviors can be efficiently predicted for novel subjects.

70 2 Results

71 2.1 CS-VAE Model Structure

72 Although existing pose estimation methods are capable enough to capture the body position of the animals
 73 in both open and contained space, tracking specific actions such as shaking and wriggling still remains a
 74 problem. However, a purely unsupervised or semi-supervised model such as a VAE or PS-VAE lacks the
 75 ability to extract meaningful and interoperable behaviors from multi-subject or social behavioral videos.
 76 One possible solution is to add another set of latent which could capture the variance across individuals
 77 and during social interactions. Instead of constraining the data points from different sessions or subjects
 78 to distinct parts of the subspace as in [2, 10], we directly constrain the latent subspace to a flexible prior
 79 distribution using a Cauchy-Schwarz regularizer as detailed in the Methods section. Ideally, this constrained
 80 subspace (CS) captures the difference between different animals in the case of a multi-subject task and the
 81 social interactions in a freely-behaving setting, while the supervised and unsupervised latents are free to
 82 capture the variables corresponding to the individual. The model structure described above is shown in Fig.
 83 1. After the input frames go through a series of convolutional layers, the resulting latent splits into three
 84 sets. The first set contains the supervised latents, which encodes the specific body position as tracked by

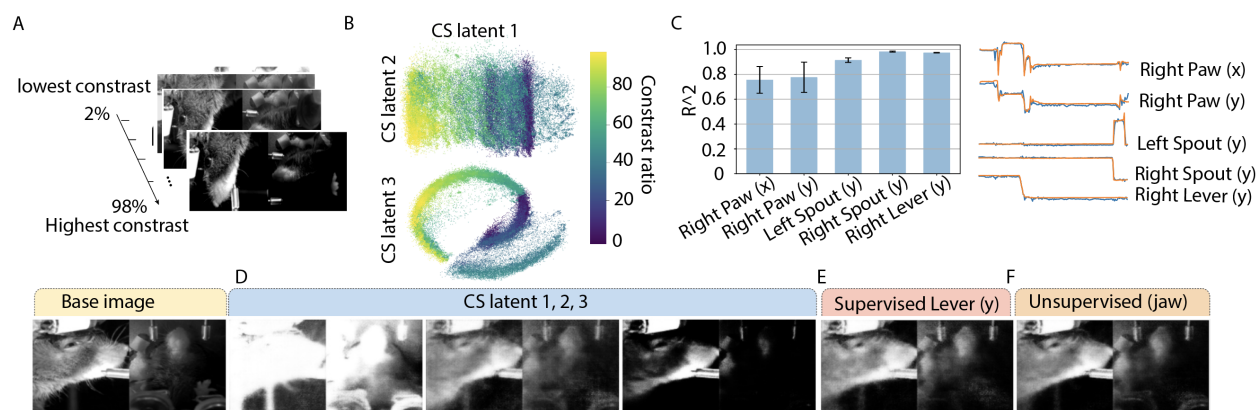


Figure 2: (A) Simulated dataset: behavioral videos from one mouse with artificially simulated differences in contrast. (B) Distribution occupied by the 3 CS latents. The constrained latents are distributed according to the pre-defined prior: a Swiss roll distribution. Different contrast ratios separate well in space. (C) Left: R^2 values for label reconstruction; Right: visualization of label reconstruction for an example trial. Latent traversals for (D) CS latents, each of which captures lower, medium, and higher contrast rate. (E) An example supervised latent captures lever movement, and (F) an example unsupervised latent which captures jaw movement.

85 supervised tracking methods such as DLC. The unsupervised latents capture the rest of the individual's
 86 behavior that are not captured by supervised latents. The CS latents capture the continuous difference across
 87 frames. The prior distribution can be changed to fit different experimental settings (and can be modeled
 88 as a discretized state space if so desired, making it close to the MSPS-VAE discussed in the Introduction).

89 2.2 Modeling Smooth Variations in a Simulated Dataset

90 We performed a simulation study on the behavioral videos of one of the mice in the ‘Multi-Subject Behavior’
 91 dataset detailed in Appendix .1. We applied a continuously varying contrast ratio throughout the trials (Fig.
 92 2A) to model smoothly varying lighting differences across the dataset. We then randomly shuffled all the
 93 trials and trained a CS-VAE model with a swiss roll as a prior distribution. Here, the R^2 for the supervised
 94 labels was 0.881 ± 0.05 (Fig. 2C), and the mean squared error (MSE) for reconstructing the entire frame was
 95 0.0067 ± 0.0003 , showing that both the images and the labels were fit well.

96 We show the CS latents recovered by the model in Fig. 2B, which follow the contrast ration distribution.
 97 We also show latent traversals in Fig. 2D-F, which demonstrate that the CS latent successfully captured
 98 the contrast changes in the frames (Fig. 2D), the supervised latent successfully captured the corresponding
 99 labeled body part (Fig. 2E), and the unsupervised latent captured parts of the individual’s body movement
 100 with a strong emphasis on the jaw (Fig. 2F). Thus, we show that smoothly varying changes in the videos are
 101 well captured by our model.

102 2.3 Modeling Multi-Subject Behavior

103 In a multi-subject behavioral task, we would like to disentangle the commonalities in behavior from the
 104 differences across subjects. Here, we test the CS-VAE on an experimental dataset with four different mice
 105 performing a two-alternative forced choice task (2AFC): head-fixed mice performed a self-initiated visual
 106 discrimination task, while the behavior was recorded from two different views (face and body). The behavioral
 107 video includes the head-fixed mice as well as experimental equipment such as the levers and the spouts. We
 108 labeled the right paw, the spouts, and the levers using DLC [3]. Neural activity in the form of widefield
 109 calcium imaging across the entire mouse dorsal cortex was simultaneously recorded with the behavior. The
 110 recording and preprocessing details are in [14, 15], and the preprocessing steps for the neural data are detailed
 111 in [15].

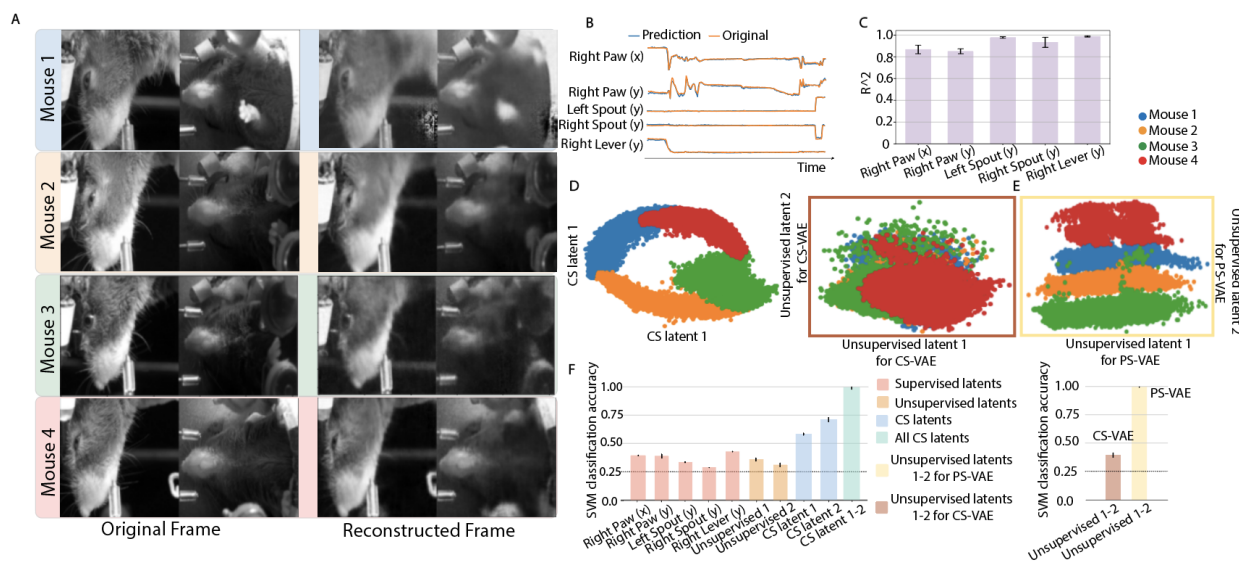


Figure 3: Modeling the behavior of four different mice. A. Image reconstruction result for an example frame from each mouse. B. Label reconstruction result for an example trial. C. R^2 value for label reconstruction for all mice. D. (Left) CS latent and (Right) unsupervised latent distributions for all mice generated using our CS-VAE model. On the left, we see that the CS latent distribution follows the pre-defined prior distribution and is well separated; on the right, we see that the unsupervised latent distribution is well overlapped across mice. E. Unsupervised latent distribution for all mice generated using the comparison PS-VAE model, where the latents from different mice are separate from each other. F. SVM classification accuracy for classifying different mice using the CS-VAE and PS-VAE latents. The unsupervised latents generated by the CS-VAE has low classification accuracy, indicating across-subject representations, and the CS latents have a classification accuracy close to one, indicating good separation.

112 **Reconstruction Accuracy** The CS-VAE model results in a mean label reconstruction accuracy $R^2 =$
 113 0.926 ± 0.02 (Fig. 3B,C), with the MSE for frame reconstruction as $0.00232 \pm 7.7 \cdot 10^{-5}$ (Fig. 3A). This was
 114 comparable to the results obtained using a PS-VAE model ($R^2 = 0.99 \pm 0.004$, $MSE = 0.13 \pm 4.5 \cdot 10^{-7}$).

115 **Disentangled Latent Space Representation** We show latent traversals for each mouse in Fig. 4, with
 116 the base image chosen separately for each mouse (videos in Supplementary Material 3). We see that, even
 117 for different mice, the supervised latent can successfully capture the corresponding labeled body part (Fig.
 118 4A). The example unsupervised latent is shown to capture parts of the jaw of each mouse (Fig. 4B), and
 119 is well-localized, comparable with the example supervised latent. The CS latent dimension encodes many
 120 different parts of the image, and has a large effect on the appearance of the mouse, effectively changing the
 121 appearance from one mouse to another, signifying that it is useful in the case of modeling mouse-specific
 122 differences (Fig. 4C). We demonstrate the abilities of the CS latent in capturing the appearance of the
 123 mouse by directly changing the CS latent from one part of subspace to another (Figure 4D). The changes in
 124 appearance along with the invariance in actions shows the intraoperability between mice by only changing
 125 the CS latents in this model (Fig. 4D).

126 Ideally, we would like to uncover common across-subject variables using the supervised and unsupervised
 127 latents subspaces, and have the individual differences across subjects be encoded in the CS latents. Thus,
 128 we expect the unsupervised latents to not be able to classify the individual well. In fact, Fig. 3D,F show
 129 that the unsupervised latents overlap well across the four mice and perform close to chance level (0.25) in
 130 a subject-classification task using SVM (details in Appendix ??). This signifies that unsupervised latents
 131 occupy the same values across all four mice and thus effectively capture across-subject behavior. In fact, we
 132 tested our latent space by choosing the same base image across the four mice, and found that the supervised
 133 and unsupervised latents from different mice can be used interchangeably to change the actions in the videos,
 134 also showing interoperability between different mice in these latent subspaces (Appendix .9).

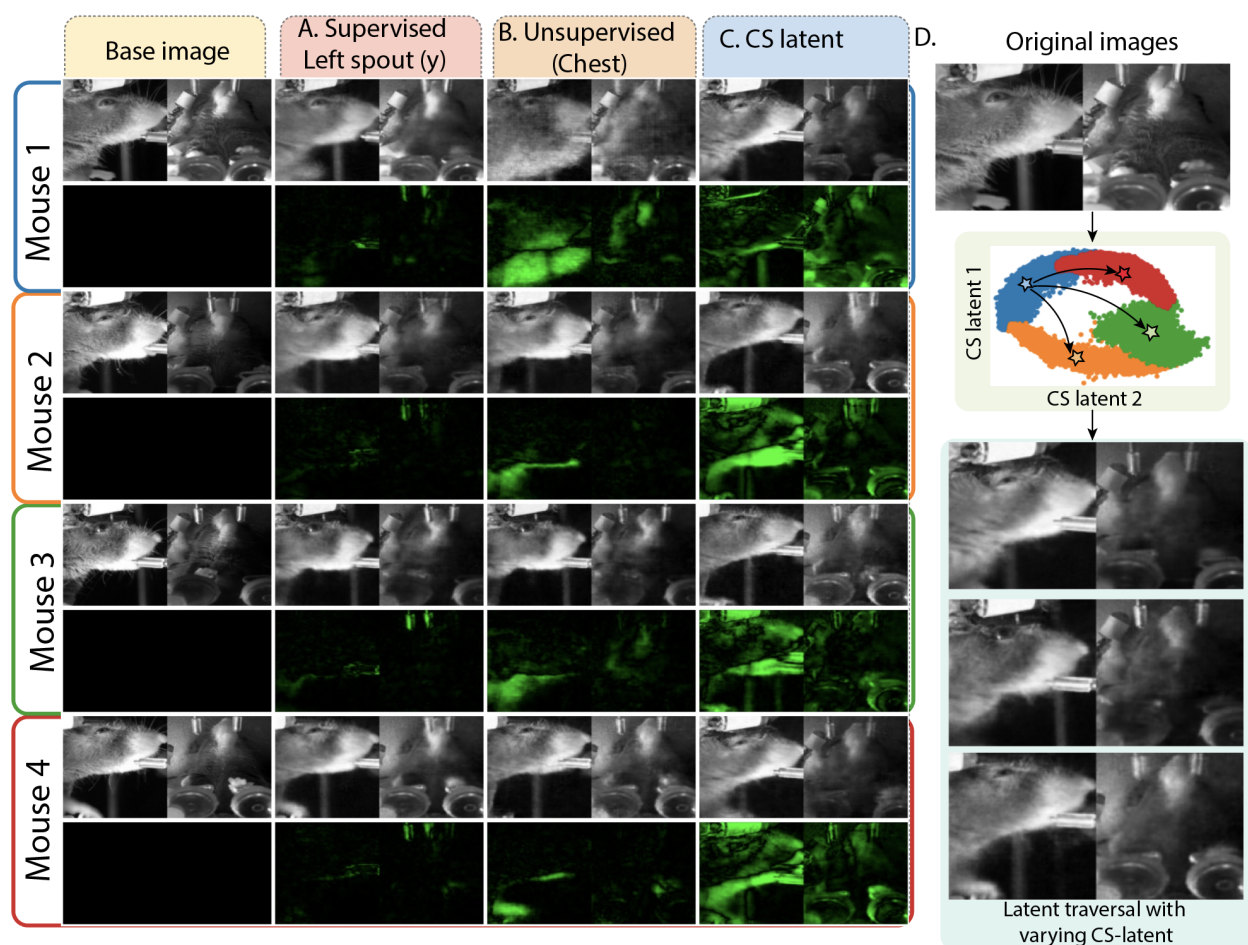


Figure 4: Latent traversals for behavioral modeling of four different mice for A. an example supervised latent that captures the left spout across all the subjects, B. an example unsupervised latent that captures the chest of the mice, and C. an example CS latent that successfully captures the mouse appearance. D. Changing the value of the CS latent in an example frame leads to a change in subject, while keeping the same action as in the example frame.

135 This is in stark contrast to the CS latents, which are well separated across mice and are able to be
 136 classified well (Fig. 3D,F); thus, they effectively encode for individual differences across subjects. Note that
 137 our method did not *a priori* know the identity of the subjects, and thus this shows that the CS latents achieve
 138 separation in an unsupervised manner. We also note that the CS latents are distributed in the shape of the
 139 chosen prior distribution (a circle). The separation in the unsupervised latent space obtained by the baseline
 140 PS-VAE shown in Fig. 3E and the latents' ability to classify different subjects (Fig. 3F) further validates the
 141 utility of CS-VAE.

142 Lastly, we trained the model while using prior distributions of different types, to understand the effect on
 143 the separability of the resulting latents. The separability was comparable across a number of different prior
 144 distributions, such as a swiss roll and a plane, signifying that the exact type of prior distribution does not
 145 play a large role.

146 **Across-Subject Motif Generation** To further show that the supervised and unsupervised latents
 147 produced by CS-VAE are interoperable between the different mice, we apply a standard SLDS model
 148 (Appendix .6) to uncover the motifs using this across-subject subspace. As seen in the ethograms (left)
 149 and the histograms (right) in Fig. 5, the SLDS using the CS-VAE latents captures common states across
 150 different subjects, indicating that the latents are well overlapped across mice. The supervised latents related

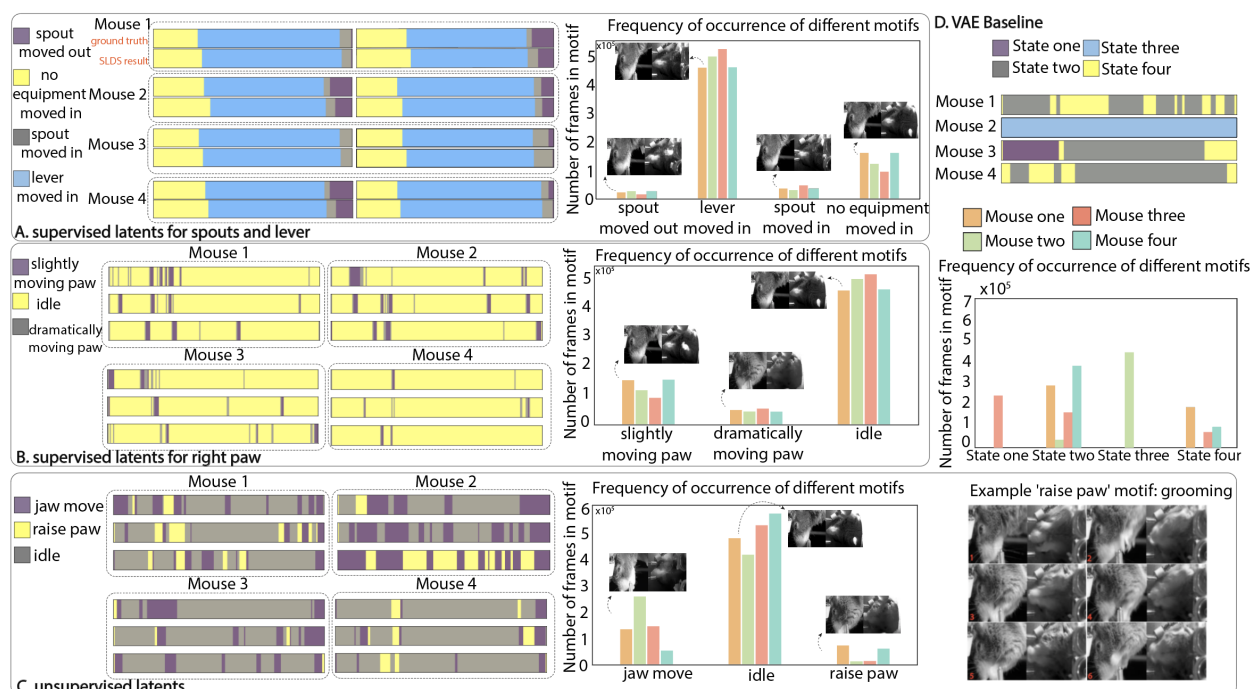


Figure 5: Motif generation for across-subject (supervised and unsupervised) behavioral latents using CS-VAE. SLDS results for CS-VAE latents: A. Supervised latents relating to equipment in the field of view. The equipment actions are similar for each trial. B. Supervised latents relating to tracked body parts. The ethograms for each trial across subjects and between subjects are very similar. The histogram indicates the number of frames occupied by each action per mouse. This further confirms the similarities between the supervised latents across subjects. C. Unsupervised latents also look similar across mice. Here, some example consecutive frames from the 'raise paw' motif are shown, which show the mouse grooming. D. As a comparison, SLDS results for the latents generated by a VAE, which failed to produce cross-subject motifs.

151 to equipment in the experiment, here the spout and lever, split the videos into four states (different colors in
 152 the ethograms in Fig. 5A), that we could independently match with ground truth obtained from sensors in
 153 these equipment. The histograms show that, as expected, these states occur with a very similar frequency
 154 across mice. We also explored the behavioral states related to the right paw. The resulting three states
 155 captured the idle vs. slightly moving vs. dramatically moving paw (Fig. 5B). The histograms show that
 156 these states also occur with a very similar frequency across mice. Videos for all these states are available in
 157 Supplementary Material 2. The inference drawn from supervised latents is directly proportional to the DLC
 158 labels. Hence, a similar conclusion can be arrived at by utilizing the DLC pose estimations. Nonetheless,
 159 the subsequent outcomes cannot be attained solely based on the poses. We extracted the behavioral states
 160 related to the unsupervised latents, which yielded 3 states related to raising of the paws (including grooming)
 161 and jaw movements (including licking) that are present in all four mice, as shown in Fig. 5C. We see that
 162 different mice have different tendencies to lick and groom, e.g., mouse 1 and 4 seem to groom more often.

163 As a baseline, we repeat this exercise on the latents of a single VAE trained to reconstruct the videos of
 164 all four mice (Fig. 5D). We see that the latents obtained by the VAE do not capture actions across subjects,
 165 and fail to cluster the same actions from different subjects into the same group.

166 **Efficient Neural Decoding via Transfer Learning** To understand the relationship between neural
 167 activity and behavior, we decoded each behavioral latent with neural data across the dorsal cortex recorded
 168 using widefield calcium imaging. The decoding results for the supervised latents were similar across the
 169 CS-VAE and the PS-VAE, but we show that the neural data was also able to capture the CS-VAE unsupervised
 170 latents well (Appendix .10).

171 Next, as a final test of interoperability of the individual latents across mice, we used a transfer learning

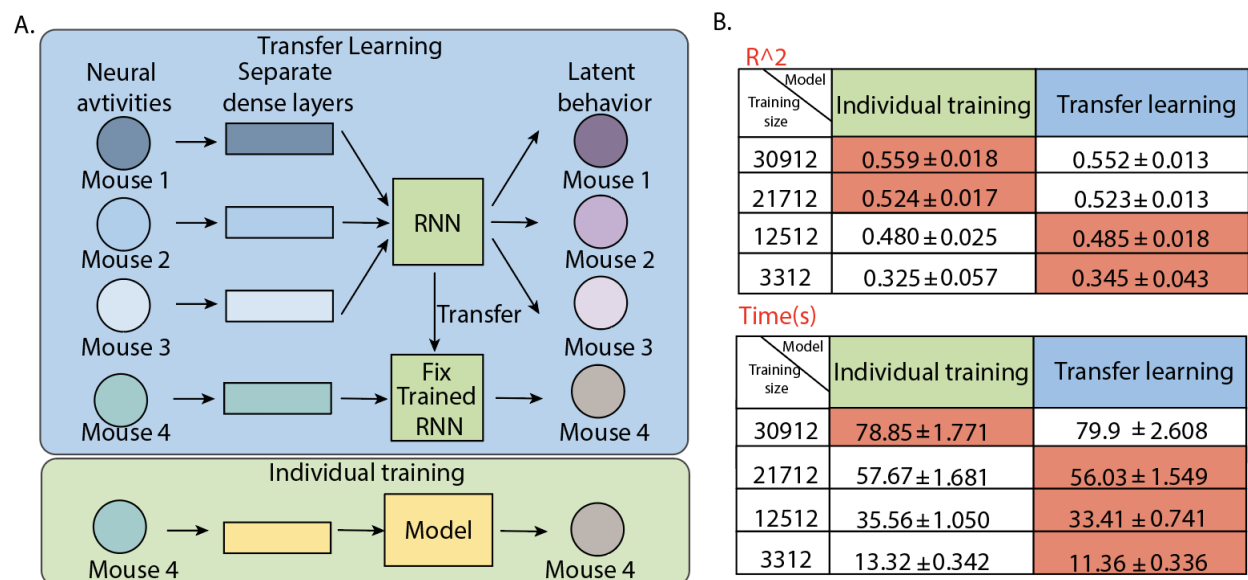


Figure 6: A. Transfer learning model framework. Each of the four mice has a specific dense layer for aligning the neural activities. After the model is trained using three mice, the across-subject Recurrent Neural Network (RNN) layer is fixed and transferred to the fourth mouse. As a comparison, we trained a novel RNN model for the fourth mouse and compared the accuracy with the transfer learning model B. R^2 and training time trade-off for individual vs. transfer learning model as the size of the training set decreases. As the training set decreases, the transfer learning has a better performance than the individually trained model with regards to both time and R^2 accuracy.

172 approach. We first trained an LSTM decoding model on 3 of the 4 mice, and then tested that model on the
 173 4th mouse while holding the LSTM weights constant but training a new dense layer leading to the LSTM (Fig.
 174 6A, details in Appendix .10). As a baseline, we compared the performance of an individual LSTM model
 175 trained only on the 4th mouse's data. We see in Fig. 6B that, as the training set of the 4th mouse becomes
 176 smaller, the transfer learning model outperforms the baseline with regards to both time and accuracy (more
 177 results and baseline comparisons in Appendix .10).

178 **Neural Correlations across Mice during Spontaneous and Task-Related Behaviors** Here
 179 we explore the neural activity correlations while the subjects perform similar spontaneous behaviors vs.
 180 task-related behaviors. Across mice, we automatically identify spontaneous behaviors such as grooming and
 181 task-related behaviors such as lever pulls. We first separate the behavior from the same motif into small
 182 segments and kept the segments that have similar means and standard deviations within and across animals
 183 as shown in Fig. 7B. Next, we explore the commonalities between the neural activity of different mice as they
 184 perform these tasks by transforming the neural activity into a common subspace, using Multidimensional
 185 Canonical Correlation Analysis (MCCA). Here, we adopt the assumptions in Safaie et al. [?] that when the
 186 animals perform the same actions, the neural latent will share similar dynamics. We employ MCCA to align
 187 the high-dimensional neural activity across multiple subjects[?]. To do this, MCCA projects the datasets
 188 onto a canonical coordinate space that maximizes correlations between them (Fig. 7 C. method details in
 189 Appendix .11). Finally, we compare the commonalities across different trials in the same subject to those
 190 across subjects for different types of behaviors. In Fig. 7D, we see that for the idle behavior, the neural
 191 correlation across mice is much lower than the correlation within the same mouse; however, this does not hold
 192 for the task-related behaviors such as lever pull and licking, or the spontaneous behaviors such as grooming.
 193 For the grooming behavior, the neural correlations within and across subjects are much higher than for the
 194 idle behaviors, and in fact, even higher than the task-related behaviors. This may be due to innate behaviors
 195 having common neural information pathways across mice, whereas learnt behaviors may display significant
 196 differences across mice. Considering the region-based differences in commonalities, the sensory areas such

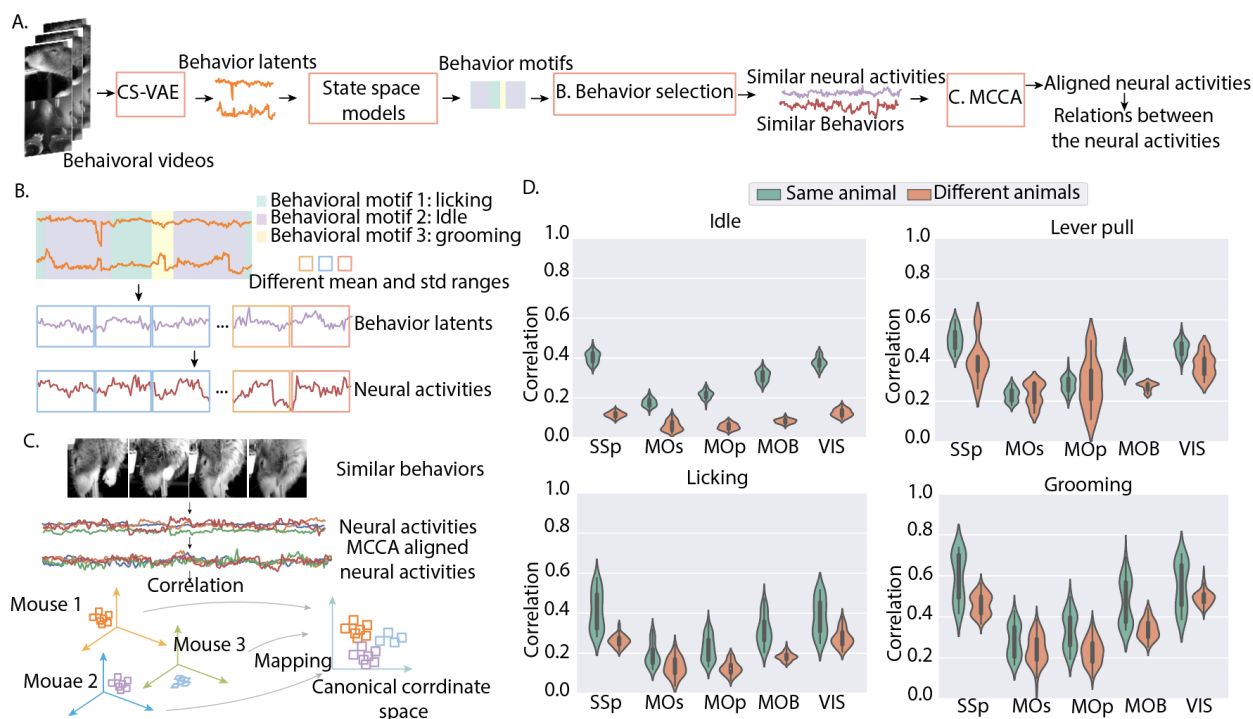


Figure 7: A. The overall workflow for comparing the neural activities for different subjects performing similar spontaneous behaviors: First, the behavioral videos are encoded into behavior latents by CS-VAE. Then, the behavior latents would be clustered into different motifs. After that, similar behaviors are grouped based on their mean and standard deviation values. We can therefore obtain the corresponding neural activities. Finally, the neural activities from different subjects are aligned using the MCCA. B. Behavior latents are cut into small fragments. Similar behavior fragments are grouped together based on their mean and standard deviation values. The corresponding neural activities are obtained based on the grouping results of the behavior. C. Neural activities are being aligned using MCCA. MCCA aligns the neural activities from different subjects by mapping them into the same feature spaces. D. Correlation score for behavioral-based aligned neural activity. The grooming behavior has higher neural correlation scores for cross-subjects than other behaviors.

197 as the visual and the somatosensory areas are much more highly correlated across mice for all behaviors as
 198 compared to motor behaviors. This may be due to the similarities in sensory feedback due to these similar
 199 behaviors but is a topic of future exploration.

200 2.4 Modeling Freely-Moving Social Behavior

201 The dataset consists of a 16 minute video of two adult novel C57BL/6J mice, a female and a male, interacting
 202 in a clean cage. Prior to the recording session the mice were briefly socially isolated for 15 minutes to increase
 203 interaction time. As preprocessing, we aligned the frame to one mouse and cropped the video (schematic
 204 in Fig. 8A; details in the Appendix .2). We tracked the nose position (x and y coordinates) of the mouse
 205 using DLC. Here, we did not include an unsupervised latent space, since the alignment and supervised labels
 206 resulted in the entire individual being explained well using the supervised latents.

207 **Reconstruction Accuracy** The CS-VAE model results in a mean label reconstruction accuracy $0.961 \pm$
 208 0.0017 (Fig. 8B), with the MSE for frame reconstruction as $1.21 \cdot 10^{-5}$ (Fig. 8B). We compared the
 209 performance of our model with the VAE and PS-VAE (Table 4), and the CS-VAE model performed better
 210 than the baseline models for both image and label reconstruction. For the VAE, we obtained the R^2 for nose
 211 position prediction by training a multi-layer perceptron (MLP) with a single hidden layer from the VAE

Table 1: Comparison of different models on the freely-moving social behavior dataset

	VAE	PS-VAE	CS-VAE
MSE for image reconstruction	$1.74 \cdot 10^{-5}$	$5.44 \cdot 10^{-5}$	$1.21 \cdot 10^{-5}$
R^2 for nose position	0.135 ± 0.013	0.894 ± 0.002	0.958 ± 0.002
R^2 for inter-individual nose-to-tail distance	0.353 ± 0.0099	0.283 ± 0.013	0.363 ± 0.0098

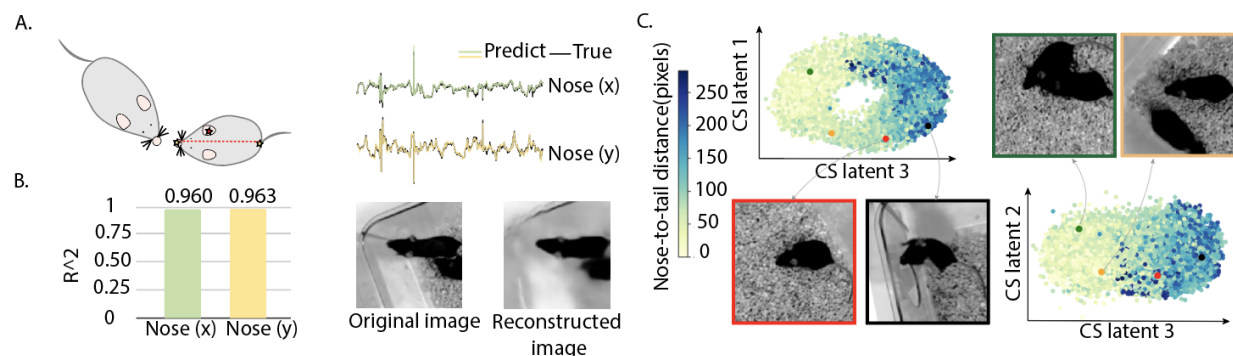


Figure 8: A. Image alignment for the social behavior data. B. Model performance on the social behavior dataset. C. Visualization of the CS latents overlaid with the nose-to-tail distance between the two interacting mice. The CS latents separates the frames that contain social interactions from those that do not.

212 latents to the nose position.

213 **Disentangled Latent Space Representation** We calculated the latent traversals for each latent as in
 214 Appendix .9. As shown in the videos in Supplementary Material 4, CS latent 1 captures the second mouse to
 215 the front of the tracked mouse, CS latent 2 captures the front and above position of the second mouse, and
 216 CS latent 3 captures the position where the second mouse is below the tracked mouse.

217 To visualize the latent space and understand the relationship to social interactions, we plot the CS latents
 218 overlaid with the nose-to-tail distance between the two mice (nose of one mouse to the tail of the other) in Fig.
 219 8C. We see that the CS latents represent the degree of social interaction very well, with a large separation
 220 between different social distances. Furthermore, we trained an MLP with a single hidden layer from different
 221 models' latents to the nose-to-tail distance, and the CS-VAE produces the highest accuracy (Table 4).

222 **Motif Generation** We applied a hidden Markov model (HMM) to the CS latents to uncover behavioral
 223 motifs. The three clusters cleanly divide the behaviors into social investigation vs. non-social behavior vs.
 224 non-social behavior with the aligned mice exploring the environment. To effectively visualize the changes
 225 in states, we show the ethogram in Fig. 9A. Videos related to these behavioral motifs are provided in
 226 Supplementary Material 5.

227 Lastly, we calculated different metrics to quantitatively evaluate the difference between each behavioral
 228 motif. The results are shown in Fig 9B, where we plot the average values for distances and angles between
 229 different key points. The lower distance between the two mice in *State a* demonstrates that the mice are
 230 close to each other in that state, pointing to social interactions. The smaller nose-to-tail distance for the
 231 aligned mouse in *State c* points to this state encoding for the 'rearing' of the mouse. The angle between the
 232 two mice further reveals the relative position between the two mice; in *State b*, the second mouse is located
 233 above the aligned mouse, while the opposite is true for *State c*. These metrics uncover the explicit differences
 234 between the different motifs that are discovered by CS-VAE.

235 3 Discussion

236 In the field of behavior modeling, there exist three major groups of methods, supervised, unsupervised, and
 237 semi-supervised. The supervised methods consist of methods such as DeepLabCut (DLC) [7], LEAP [6],

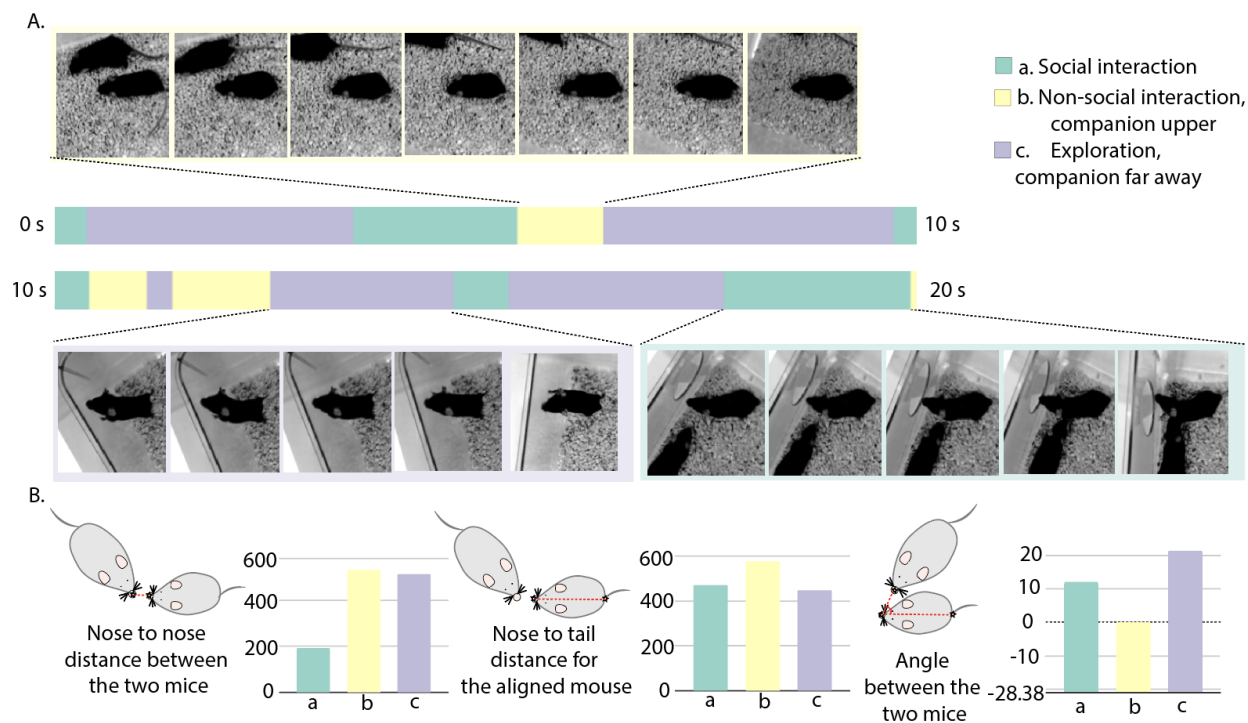


Figure 9: A. Ethogram for the animals' behavior recovered using hidden Markov models (HMM) applied to the CS latents. B. Different metrics for analysing the behavioral motifs. Here, the three motifs are *a*. social interaction; *b*. non-social interaction with the companion on the upper side of the aligned mouse; *c*. non-social interaction (the aligned mouse exploring the environment with its companion far away). These metrics show the quantitative differences between the different motifs.

238 AlphaTracker [5], amongst others. Although these methods capture the positions of the subjects, they lack the
 239 ability to model smaller movements and unlabeled behavior, and necessitate tedious labeling. On the other
 240 hand, unsupervised methods such as MoSeq [9] and Behavenet [8] lack the ability to produce interpretable
 241 behavioral latents. While some semi-supervised methods, for instance, MSPS-VAE [2] and DBE [10], succeed
 242 in producing interpretable latents and modeling behavior across subjects, they need significant human input,
 243 and lack the ability to model freely-moving animals' behavior. Here, we introduce a constrained generative
 244 network called CS-VAE that effectively addresses major challenges in behavioral modeling- disentangling
 245 multiple subjects and representing social behaviors.

246 For multi-subject behavioral modeling, the behavioral latents successfully separates the common activities
 247 across animals from the differences across animals. This behavioral generality is highlighted by the across-
 248 subject behavioral motifs generated by standard methods, and a higher accuracy while applying transfer
 249 learning for the neural decoding task. Furthermore, the SVM classification accuracy approaches 100%, which
 250 also indicates that the constrained-subspace latents well separate the differences between the subjects. In
 251 the social behavioral task, the constrained latents well capture the presence of social investigations, the
 252 environmental exploration, and the relative locations of the two individuals in the behavioral motifs. While
 253 our methods succeed in effectively modeling social behavior, it remains a challenge to separate out different
 254 kinds of social investigations in an unsupervised manner.

255 The constrained latents encode smoothly and discretely varying differences in behavioral videos. As seen
 256 in this work, in the across-subject scenario, the constrained latents encode the appearance of the different
 257 subjects, while in freely-moving scenario, the constrained latents capture social investigation between the
 258 subjects. The flexibility of this regularization thus gives it the ability to be fit in different conditions. Future
 259 directions include building an end-to-end structure that can captures behavioral motifs in a unsupervised way.

260 4 Methods

261 **Regularization of Constrained Subspace** We use the Cauchy-Schwarz divergence to regularize our
 262 constrained subspace using a chosen prior distribution. The Cauchy-Schwarz divergence $D_{CS}(p_1, p_2)$ between
 263 distributions $p_1(x)$ and $p_2(x)$ is given by:

$$D_{CS}(p_1, p_2) = -\log \frac{\int p_1(x)p_2(x)dx}{\sqrt{\int p_1^2(x)dx \int p_2^2(x)dx}} \quad (1)$$

264 $D_{CS}(p_1, p_2)$ equals zero if and only if the two distributions $p_1(x)$ and $p_2(x)$ are the same. By applying the
 265 Parzen window estimation technique to $p_1(x)$ and $p_2(x)$, we get the entropy form of the Equation [11]:

$$\hat{H}(p_1) = -\log(V(p_1)) = -\log \left(\sum_i^N \sum_j^N G_{\sqrt{2}\sigma}(p_{1i} - p_{1j})/N^2 \right) \quad (2)$$

$$\hat{H}(p_1, p_2) = -\log(V(p_1, p_2)) = -\log \left(\sum_i^{N_1} \sum_j^{N_2} G_{\sqrt{2}\sigma}(p_{1i} - p_{2j})/(N_1 N_2) \right) \quad (3)$$

266 Here, p_{1i} represents the i th sample from the distribution p_1 , i.e., $p_1(x_i)$. $-\log(V(p_1))$ and $-\log(V(p_2))$ are
 267 the estimated quadratic entropy of $p_1(x)$ and $p_2(x)$, respectively, while $-\log(V(p_1, p_2))$ is the estimated
 268 cross-entropy of $p_1(x)$ and $p_2(x)$. G is the kernel applied to the input distribution; here it is chosen to be
 269 Gaussian. N , N_1 , and N_2 are the number of samples being input into the model while σ is the kernel size.
 270 The choice of the kernel size depends on the dataset itself; generally, the kernel size should be greater than
 271 the number of the groups in the data. Equation (1) can be expressed as:

$$\mathcal{L}_{CS} := D_{CS}(p_1, p_2) = \log \frac{V(p_1)V(p_2)}{V^2(p_1, p_2)} \quad (4)$$

272 Here, $p_1(x)$ represents the distribution of our CS latent space, and $p_2(x)$ the chosen prior distribution. In
 273 Equation (4), minimizing $V(p_1)$ would result in the spreading out of $p_1(x)$, while maximizing $V(p_1, p_2)$ would
 274 make the samples in both distributions closer together [11]. Thus, we minimize this term in the objective
 275 function while training the model. However, it may be necessary to stop at an appropriate value, since overly
 276 spreading out $p_1(x)$ may lead to the separation of the samples from the same groups, while making p_1 and p_2
 277 excessively close may cause mixtures of data points across groups.

278 In short, the Cauchy-Schwarz divergence measures the distance between p_1 and p_2 . In our work, we adopt
 279 a variety of distributions as a prior distribution $p_2(x)$, and we aim to project the constrained subspace latents
 280 onto the prior distribution (see Fig. 1).

281 **Optimization** The loss for the CS-VAE derives from that for the PS-VAE, and is given by:

$$\mathcal{L}_{CS-VAE} = \mathcal{L}_{frames} + \alpha \mathcal{L}_{label} - \mathcal{L}_{KL-s} - \mathcal{L}_{ICMI} - \beta \mathcal{L}_{TC} - \mathcal{L}_{DWKL} + \gamma \mathcal{L}_{CS} \quad (5)$$

282 Here, the terms \mathcal{L}_{frames} and \mathcal{L}_{label} represent the reconstruction loss of the frames and the labels, respectively.
 283 The \mathcal{L}_{KL-s} represents the KL-divergence loss for the supervised latents while \mathcal{L}_{ICMI} , \mathcal{L}_{TC} , and \mathcal{L}_{DWKL}
 284 form the decomposed version of the KL loss for the unsupervised latents. Lastly, the \mathcal{L}_{CS} represents the CS-
 285 divergence loss on our constrained latents. α is introduced to control the reconstruction quality of the labels, β
 286 is adopted to assist the model in producing independent unsupervised latents, and γ is implemented to control
 287 the variability in the constrained latent space for better separation. The detailed explanations and derivations
 288 for each term in the objective function are in Appendix .3. Furthermore, the loss terms in Equation (5) can be
 289 modified to fit various conditions. For a freely-behaving social task, the background for one individual in the
 290 container could be the edge of the container as well as the rest of the individuals in the container. The choice of
 291 hyperparameters and the loss curves through the training process is shown in Appendix .5 and .7, respectively.

292 **Visualization of the latent space** To test how the image varies with a change in the latent, one frame
 293 from the trials is randomly chosen as the ‘base image’, and the effect of varying a specific latent at a time
 294 is visualized and quantified. This is known as the ‘latent traversal’ [2]. First, for each latent variable, we

295 find out the maximum value that it occupies across a set of randomly selected trials. We then change that
296 specific latent to achieve its maximum value, and this new set of latents forms the input to the decoder. We
297 obtain the corresponding output from the decoder as the ‘latent traversal’ image. Finally, we visualize the
298 difference between the ‘latent traversal’ image and the base image. The above steps are performed for each
299 latent individually. In videos containing latent traversals (Supplementary Material), we change the latent’s
300 value from its minimum to its maximum across all trials, and input all the corresponding set of latents into
301 the decoder to produce a video.

302 **Behavioral Motif Generation** Clustering methods such as Hidden Markov Models (HMM) and switching
303 linear dynamical systems (SLDS) have been applied in the past to split complex behavioral data into simpler
304 discrete segments [16] (see Appendix .6 for details). We use these approaches to analyze motifs from our
305 latent space, and directly input the latent variables into these models. In the case of multi-subject datasets,
306 our goal is to capture the variance in behavior in a common way in the across-subject latents, i.e., recover
307 the same behavioral motifs in subjects performing the same task. In the case of freely-moving behavior, our
308 goal is to capture motifs related to social behavior.

309 **Efficient Neural Decoding** Decoding neural activity to predict behavior is very useful in the under-
310 standing of brain-behavior relationships, as well as in brain-machine interface tasks. However, models to
311 predict high-dimensional behavior using large-scale neural activity can be computationally expensive, and
312 require a large amount of data to fit. In a task with multiple subjects, we can utilize the similarities in
313 brain-behavior relationships to efficiently train models on novel subjects using concepts in transfer learning.
314 Here, we represent across-subject behavior in a unified manner and train an across-subject neural decoder.
315 Armed with this across-subject decoder, we show the decoding power on a novel subject with varying amounts
316 of available data, such that it can be used in a low-data regime. The implementational details for this transfer
317 learning approach can be found in Appendix .10.

318 **Behavior election for innate behaviors studying** While the behavioral features extracted from
319 the previous sections are successful in capturing similar spontaneous behaviors across various animals, the
320 behavioral patterns within the same motifs can exhibit substantial variation. For instance, in the case of
321 the raising paw motif, continuous movement of the paws could be indicative of either grooming or other
322 complex behaviors. To overcome this challenge, we divided the behaviors belonging to the same motif into
323 smaller segments and calculated the corresponding mean and standard deviation of the behavioral latents.
324 Subsequently, we compared these values and retained the segments that exhibited similar mean and standard
325 deviations both within and across animals, as illustrated in Fig. 7B. These steps were repeated for all the
326 behavioral motifs examined in this study.

327 In addition to the spontaneous behaviors discussed above, we also selected an ‘idle’ behavior that captured the
328 mouse’s inactivity and a task-related behavior, namely the ‘lever pull’ behavior, which signaled the initiation
329 of each task.

330 References

- 331 [1] Krakauer, J. W., Ghazanfar, A. A., Gomez-Marin, A., MacIver, M. A. Poeppel, D. Neuroscience needs
332 behavior: Correcting a reductionist bias. *Neuron* 93, 480–490 (2017).
- 333 [2] Whiteway, M. R. et al. Partitioning variability in animal behavioral videos using semi-supervised variational
334 autoencoders. *bioRxiv* (2021).
- 335 [3] Mathis, A. et al. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning.
336 *Nature Neuroscience* 21, 1281–1289 (2018).
- 337 [4] Pereira, T. et al. Fast animal pose estimation using deep neural networks. *bioRxiv* (2018).
- 338 [5] Chen, Z. et al. Alphatracker: A multi-animal tracking and behavioral analysis tool. *bioRxiv* (2020).
- 339 [6] Pereira, T. D. et al. Publisher correction: Slep: A deep learning system for multi-animal pose tracking.
340 *Nat Methods* (2022).
- 341 [7] Lauer, J. et al. Multi-animal pose estimation and tracking with deeplabcut. *bioRxiv* (2021).

- 342 [8] Batty, E. et al. Behavenet: nonlinear embedding and bayesian neural decoding of behavioral videos. In
343 Wallach, H. et al. (eds.) Advances in Neural Information Processing Systems, vol. 32 (Curran Associates,
344 Inc., 2019).
- 345 [9] Wiltschko, A. B. et al. Revealing the structure of pharmacobehavioral space through motion sequencing.
346 Nature neuroscience 23, 1433 –1443 (2020).
- 347 [10] Shi, C. et al. Learning disentangled behavior embeddings. In NeurIPS (2021).
- 348 [11] Santana, E., Emigh, M. Principe, J. Information theoretic-learning auto-encoder (2016).
- 349 [12] Tran, L., Pantic, M. Deisenroth, M. P. Cauchy-schwarz regularized autoencoder (2021). 2101.02149.
- 350 [13] Wiltschko, A. et al. Mapping sub-second structure in mouse behavior. Neuron 88, 1121–1135 (2015).
- 351 [14] Musall, S., Kaufman, M. T., Juavinett, A. L., Gluf, S. Churchland, A. K. Single-trial neural dynamics
352 are dominated by richly varied movements. Nature neuroscience 22, 1677 – 1686 (2019).
- 353 [15] Saxena, S. et al. Localized semi-nonnegative matrix factorization (locanmf) of widefield calcium imaging
354 data. PLOS Computational Biology 16, 1–28 (2020).
- 355 [16] Linderman, S. et al. Bayesian Learning and Inference in Re- current Switching Linear Dynamical Systems.
356 In Singh, A. Zhu, J. (eds.) Proceedings of the 20th International Conference on Artificial Intelligence and
357 Statistics, vol. 54 of Proceedings of Machine Learning Research, 914–922 (PMLR, 2017).

358 5 Appendix

359 .1 Experimental Methods and Preprocessing for the Multi-Subject Dataset

360 In our work, we employed a subset of the behavioral dataset detailed in Musall et al., 2019 [14]. Briefly, the
361 task entailed pressing a lever to initiate the task, after which a visual stimulus was displayed towards the left
362 or the right. After a delay period, the spouts come forward, at which time the mouse makes its decision by
363 licking the spout corresponding to the direction of the visual stimulus (left or right). Finally, the mice receive
364 a juice reward if they choose correctly.

365 We tested the CS-VAE on the behavioral data for the four mice performing a visual task and randomly
366 chose 388 trials per mouse each of the trials has a 189 number of frames. Each frame was pre-processed and
367 resized to have both the length and width being 128. One example trial for each mouse can be found in
368 Supplementary Material 1.

369 Before inputting the data into the model, we sorted the trials by the amount of variance in the images,
370 and shuffled the first half (high variance) and the second half (low variance) of the dataset separately. This
371 was done to speed up training by training the model on high-variance trials first. We tested our model by
372 randomly choosing 4 trials from all trials for each mouse 5 times. The same procedure was applied when
373 training the model on the simulation dataset, i.e., the doctored data for one subject.

374 .2 Experimental Methods and Preprocessing for the Freely-Moving Social Be- 375 havior Dataset

376 The dataset consists of a 16-minute video of two adult novel C57BL/6J mice, a female and a male, interacting
377 in a clean cage. Prior to the recording session, the mice were briefly socially isolated for 15 minutes to increase
378 interaction time. This dataset was collected by one of the authors. The original data has 24917 number of
379 frames with length and width being 1920 and 1080, respectively. The example fraction of the video can be
380 found in Supplementary Material 6.

381 The nose, ears, and tail base of each mouse were manually annotated using AlphaTracker. We kept 19659
382 number of frames that have the labels for preprocessing and training. We perform several preprocessing steps
383 to align and crop the video as well as the labels based on one of the two mice (Mouse 1, female). All of the
384 preprocessing steps were based on the AlphaTracker labels. For each frame, we first rotate it to ensure that
385 the nose and tailbase for Mouse 1 are on the same horizontal line, with the central point for rotation as the
386 left ear. Next, we aligned the frame such that the left ear of Mouse 1 was at the same location across all
387 frames. Finally, we resize the frame to be 128×128 and consequently the AlphaTracker labels. For this
388 dataset, since there was a relatively low number of frames, we obtained the CS-VAE MSE and label R^2 for
389 the entire dataset.

390 .3 Methodological details of the Partitioned Subspace VAE

391 The Partitioned Subspace VAE (PS-VAE) was introduced in [2], and we borrow the notation used in that
392 paper when detailing the CS-VAE. Thus, we include here a full description of the model.

393 First of all, we define the input frame as x , and the corresponding pose estimation tracking label as y .
394 The reconstructed variables are termed \hat{x} and \hat{y} , respectively. The supervised latent space is denoted as z_s ,
395 unsupervised latent as z_u , and the background latent as z_b . In a VAE model, we would like to minimize the
396 distance, typically the KL divergence, between the posterior distribution of the latent variables $p(z|x)$ and a
397 chosen distribution $q(z|x)$. However, since $p(z|x)$ is an unknown distribution, the Evidence Lower Bound
398 (ELBO) is introduced as an alternative method to reduce the KL divergence:

$$399 \mathcal{L}'_{ELBO} = \mathbb{E}_{q(z|x)}[\log(p(x|z))] - KL[q(z|x)||p(z)] \quad (6)$$

400 Following [2], if we have a finite dataset $\{x_n\}_{n=1}^N$, and we treat n as a random variable with a uniform
401 distribution $p(n)$ while defining $q(z|n) := q(z_n|x_n)$, we can rewrite the ELBO as:

$$\mathcal{L}_{ELBO} = \mathbb{E}_{p(n)}[\mathbb{E}_{q(z|n)}[\log(p(x|z))] - \mathbb{E}_{p(n)}[KL[q(z|n)||p(z)]] \quad (7)$$

402 We define the loss over frames \mathcal{L}_{frames} as the first of the two terms above. In the PS-VAE model, there are
 403 two inputs: frames x and labels y . Therefore, in Equation (7), instead of writing the input likelihood as
 404 $p(x|z)$, we can now write it as $p(x, y|z)$. A simplifying assumption is made that x and y are conditionally
 405 independent given z , and thus we can directly write $\mathcal{L}_{frames+labels}$ as $\mathcal{L}_{frames} + \mathcal{L}_{labels}$, where \mathcal{L}_{labels} is
 406 calculated by replacing x with y in \mathcal{L}_{frames} .

After assuming the prior $p(z)$ has a factorized form: $p(z) = \prod_i p(z_i)$, the KL term \mathcal{L}_{KL} can be split as the addition of ℓ_{KL-s} and ℓ_{KL-u} , i.e., the KL terms for the supervised and unsupervised latents, respectively. We decompose the KL term for the unsupervised latent as the following [2].

$$\begin{aligned} \mathcal{L}_{KL-u} &= \mathcal{L}_{ICMI} + \mathcal{L}_{TC} + \mathcal{L}_{DWKL} \\ &= KL[q(z_u, n) || q(z_u)p(n)] + KL[q(z_u) || \prod_j (z_{u,j})] + KL[q(z_{u,j}) || \prod_j (z_{u,j})] \quad (8) \end{aligned}$$

407 where j represents the latent dimension, \mathcal{L}_{ICMI} is the index-code mutual information, which measures how
 408 well the latent encodes the corresponding input data. The term TC is short for total correlation, which
 409 measures the interdependency of each latent dimension. The third term, \mathcal{L}_{DWKL} is the dimension-wise KL,
 410 which calculates the KL divergence for each dimension individually. Finally, the resulting subspace is forced
 411 to be orthogonal by applying orthogonal weights across all the different latents.

412 The authors in [2] introduce an extension to PS-VAE for modeling multi-session data. The Multi-Session
 413 PS-VAE (MS-PS-VAE) can only work with a labeled set of discrete sessions, as described in the Introduction.
 414 The images from each session are labeled, and the session-specific latents are enforced to be static over
 415 time, thus capturing the image-related details. To enforce the background latents to be static over time in a
 416 particular session, and to maximize the difference in the background latents across different sessions, the
 417 triplet loss is introduced in MS-PS-VAE. As described in the Introduction, this loss term artificially places
 418 the latents from the same session together while separating the latents from different sessions. The triplet
 419 loss is computed as the following.

$$\mathcal{L}_{triplet} = \max\{d(a, p) - d(a, n) + m, 0\} \quad (9)$$

420 Here, a is the anchor point, p is the positive point, n is the negative point, and m is a margin. The function
 421 pulls the point p towards point a , and pushes the point n away from point a . While training, the data from
 422 multiple sessions is included in each mini-batch. The data from each session is split in three, and each third
 423 from the same session acts as an anchor and positive point, while the data from another session acts as a
 424 negative point. Practically, this requires as many sessions as possible in the same mini-batch during the
 425 training for accurate results. As the number of sessions increases, this method becomes computationally
 426 intractable, and may lead to unsatisfactory reconstruction results. Moreover, this loss does not allow for
 427 varying backgrounds across any one session.

428 In the MS-PS-VAE model, the triplet loss was applied as a supervised manner to pull the data from the
 429 same subject being closer while pushing the different subjects away from each other. This method is only
 430 useful when the number of sessions is known, and is not applicable in an open-field setting, for example while
 431 modeling freely-moving social behavior as in this manuscript.

432 Therefore, in this manuscript, we introduce a regularization term that can automatically separate different
 433 subjects in the background latent space without specifying the number of sessions or labeling each frame as
 434 belonging to a specific session.

435 .4 Model Architecture and Training

436 Our computational experiments were carried out using TensorFlow and Keras. The image decoder we use
 437 is symmetric to the encoder, with both of them containing 14 convolution layers. We applied the Adam
 438 optimizer with learning rate as 10^{-4} . For the multi-subject dataset, we fixed our batch size to be 256 and
 439 trained for 50 epochs. For the freely-moving social behavior dataset, we trained for 500 epochs with batch
 440 size 128.

441 .5 Choice of Hyperparameters

442 In the multi-subject dataset, four coefficients need to be decided for the objective function as indicated above:
 443 $\{\alpha, \beta, \sigma, \gamma\}$. There is a balance between the choice of β and γ : properly choosing the values could separate
 444 the latent in the unsupervised space and the latents in both unsupervised and background space as well. A
 445 large separation of the background latent may potentially lead to unsatisfactory reconstruction results. The
 446 choice of kernel size σ depends on the dataset, and should be larger than the number of distinct groups in
 447 our dataset; since in our current experiments, we have at most four groups, we set $\sigma = 15$. Moreover, we set
 448 α to 1000, β to 5, γ to 500. We set the dimensionality of the supervised latent space equal to the number of
 449 tracked video parts, which is 5 in our case. We set the dimensionality of the unsupervised latent space as 2,
 450 while that of the background latent space as 2.

451 In the social behavior task, we track the nose location as the supervised latent, since the other labels do
 452 not have a high variance (due to the alignment process). Additionally, we do not need any unsupervised
 453 latents to explain the individual’s behavior. The CS latent in this setting has 3 dimensions. Here, α is 1200,
 454 γ is 200, and the kernel size is 20.

455 The hyperparameters chosen for all three datasets are shown in Tables 2 and 3.

Table 2: Hyperparameter for different dataset

Dataset	α	β	σ	γ
Various contrast	1000	5	5	500
Multi-subject	1000	5	15	500
Social behavior	1200	N/A	20	200

Table 3: Latent dimensions and the prior distribution for different dataset

Dataset	supervised	unsupervised	constrained	prior distribution
Various contrast	5	2	3	Swiss roll
Multi-subject	5	2	2	circle
Social behavior	2	0	3	hollow cylinder

458 .6 Motif Generation

459 A switching linear dynamical system (SLDS) consists of discrete latent state $z_t \in \{1, 2, \dots, K\}$, continuous
 460 latent state $x_t \in \mathbb{R}^M$, and the observation state $y_t \in \mathbb{R}^N$. Here, $t = 1, 2, 3, \dots, T$ is the time step, T is the
 461 length of the input signal; K is the number of discrete states; M is the number of latent dimensions; N is the
 462 observation dimensions. The discrete latent state z_t follows the Markovian dynamics with the state transition
 463 matrix expressed as:

$$Q_{i,j} = P(z_t = j | z_{t-1} = i) \quad (10)$$

464 The continuous latent state x_t has the following linear dynamical relations that determined by z_t .

$$x_{t+1} = A_{z_{t+1}}x_t + V_{z_{t+1}}u_t + b_{z_{t+1}} + w_t \quad (11)$$

465 Here, $A_{z_{t+1}}$ is the dynamic matrix at state z_{t+1} ; u_t is the input at time t, with $V_{z_{t+1}}$ being the control
 466 matrix; $b_{z_{t+1}}$ is the offset vector and w_t being the noise which is generally the zero mean Gaussian. Here, our
 467 observation model is in Gaussian case; therefore, the observation y_t is expressed as:

$$y_t = C_{z_t}x_t + F_{z_t}u_t + d_{z_t} + v_t \quad (12)$$

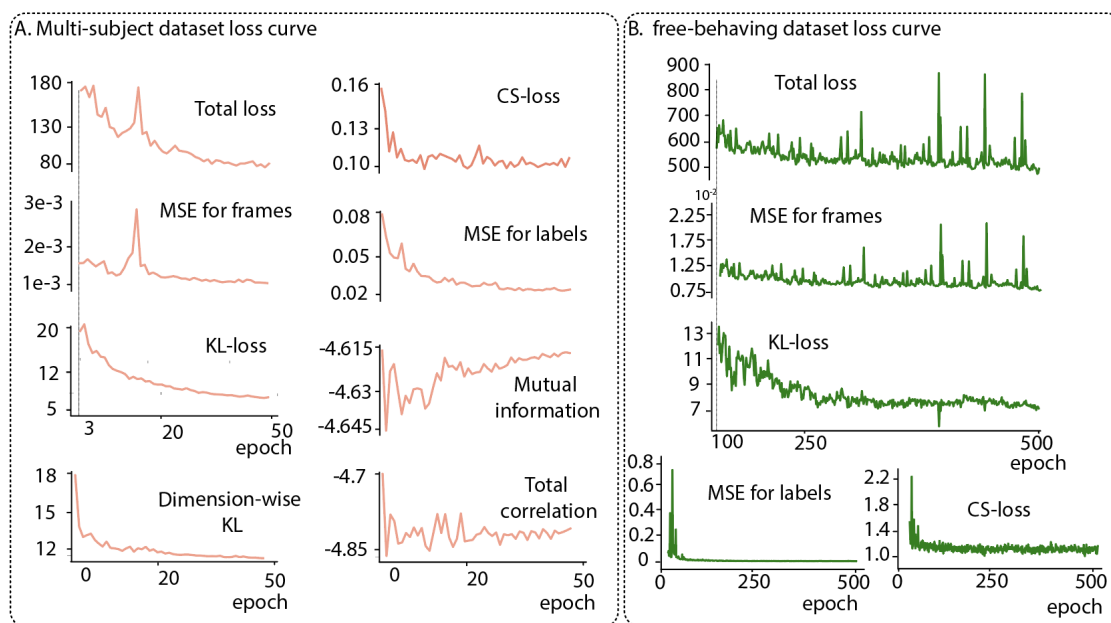
468 Here, C_{z_t} is the measurement matrix at state z_t ; F_{z_t} is the feedthrough matrix which directly feed the
 469 input into the observation; d_{z_t} is the offset vector and v_t is the noise. Here the update was accomplished by
 470 the Expectation-Maximization(EM) algorithm. In the E-step, the model updates the hyperparameters. In
 471 the M-step, the log-likelihood in Eq.12 is being maximized.

472 To implement the SLDS, we adopted the open source software from Linderman et al.[16]. We fit the
 473 SLDS using different latent dimensions, where the observation dimension was the order of latent dimension

474 and the number of states was determined by visualizing the videos. We use SLDS's to model the motifs in
 475 the multi-subject dataset since the behaviors are well separated using their dynamics. We use K-means to
 476 model the motifs in the freely-moving social behavior dataset since the behaviors are well separated directly
 477 in state space. An autoregressive HMM (a simpler model than an SLDS) applied to the CS latents in the
 478 social behavior dataset leads to similar results as the K-means.

479 .7 Loss Curves

480 We show the learning curve for each loss term for both dataset to precisely quantify the model, in Fig. 10.
 481 For the multi-subject dataset (Fig. 10A), for the unsupervised latents, the final loss for dimension-wise KL,
 482 total correlation, and the mutual information are 11.7, -4.8 , and -4.6 , respectively. The final KL loss for
 483 the supervised latents is 5.06 and the final CSD loss for the CS latents is 0.1. For the free behaving dataset,
 484 the loss curves for each loss term are shown in Fig. 10B. By the end of the training process, the KL loss for
 485 the supervised latents is 7.01 and the CSD loss for the CS latents is 1.15.



486

Figure 10: Loss curve for A. training the multi-subject dataset B. training the freely behaving dataset with the specified hyperparameters as in Tables 1 and 2.

487 .8 SVM

488 To further quantify the separation of the latents between different subjects, we applied a supervised classifica-
 489 tion method to decode the identity of the subject using each latent.

490 After randomly shuffling all the latents, we split all the trials into training trials and test trials, with
 491 each mouse having 368 trials in the training set and 20 trials in the test set, and repeated this 5 times with
 492 different random seeds.

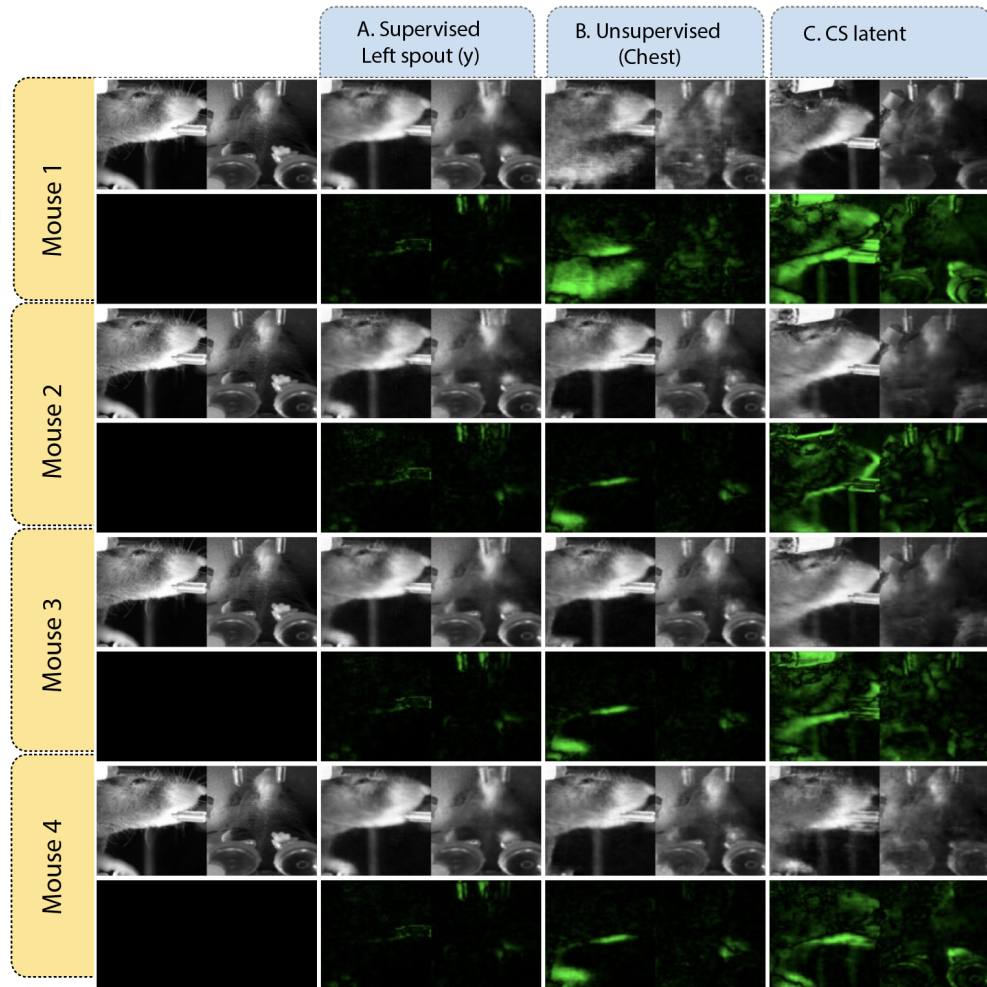
493 .9 Latent traversal

494 For the multi-subject dataset, we tested the latent traversal with the same base image to validate the results,
 495 shown in Figure 11. Here, we randomly chose a frame from a mouse and changed each individual latent
 496 within different ranges as detailed in the Methods. For example, in Figure 11, the first row contains the
 497 output when the corresponding latent is changed to take on the maximum value from the range of Mouse 1.
 498 Similar to the figures in the main text, the upper images are the latent traversal images while the lower ones
 499 are the difference between the upper and original images. We see that the base image from Mouse 3 can be

500 flexibly changed to produce a different mouse when changing the CS latent. Moreover, when changing the
501 supervised and unsupervised latents for the different mice, Mouse 3 seems to be flexibly changing with these
502 latents from different mice.

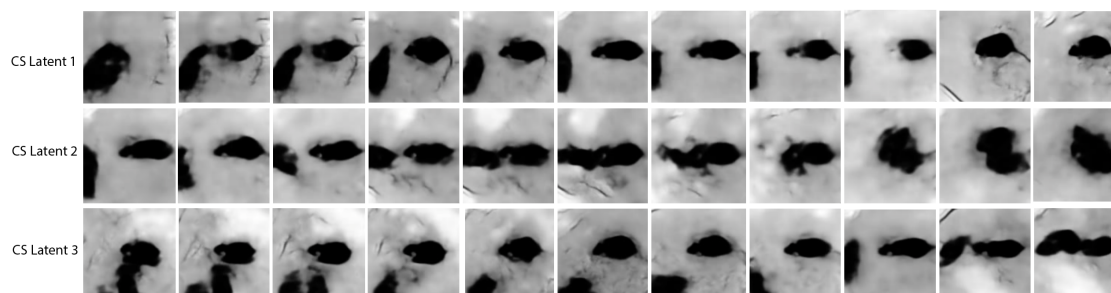
503 To better visualize the specialization of each latent, we generated the latent traversal videos for each
504 latent with different base images. For different mice, we, first of all, find the maximum and the minimum
505 value for the specific latent. Then, change the latent within that range with 0.5 per step. Finally, concatenate
506 all the latent traversal images into videos. The videos can be found in Supplementary Material 3.

507 We performed a similar visualization on the freely-moving social behavior dataset for the CS latents. The
508 latent traversal videos can be found in Supplementary Material 4, and some clips from the videos are shown
509 in Fig 12.



510

Figure 11: Latent traversals for the multi-subject dataset for the four mice with the same base image A. an example supervised latent, B. an example unsupervised latent, and C. an example CS latent. We see that the same base image (Mouse 3) is transformed into a different mouse each time when changing the CS latent.

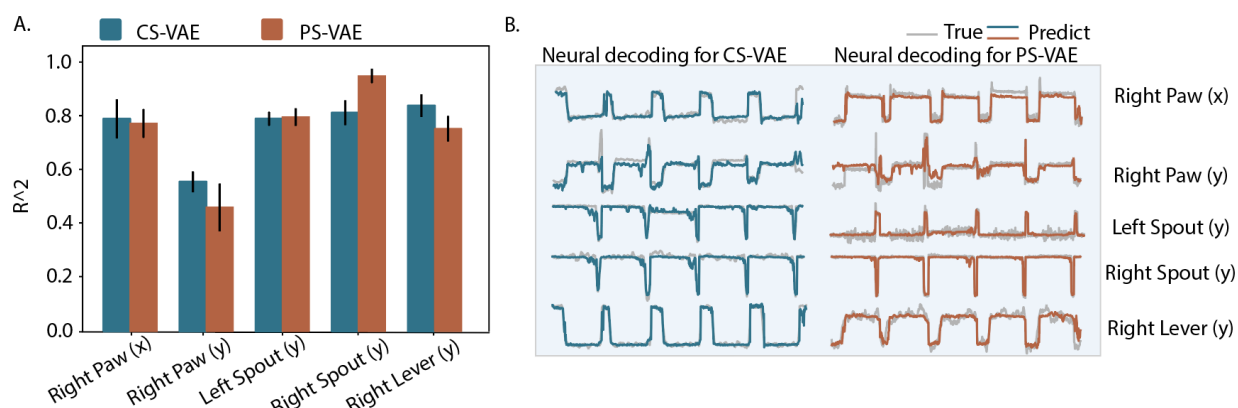


511

Figure 12: Latent traversals on the CS-latents for the freely-moving social behavior dataset. We see that the latents all encode for social interactions between the two mice.

512 .10 Neural decoding models

513 The trials were first shuffled and then split into training and testing. Next, we employed the CS-VAE
 514 generated latent representations, and choose one example subject to decode the behavior at time t using the
 515 neural activity recorded between $t - 0.15s$ and t . We applied four types of models to compare the performance.
 516 A linear model which directly maps the neural activities into the behavior. A multilayer perceptron (MLP)
 517 with three dense layers to train the decoder. We used the Adam optimizer with learning rate decay from 0.1
 518 with 0.3 decay rate for every 5 step. The batch size was fixed to be 150 and trained for 200 epochs. A LSTM
 519 model, which begin with a dense layer followed by a LSTM layer with a drop-out rate being 0.5 and another
 520 dense layer at the end. We applied the same training strategy as in MLP model.



521

Figure 13: Neural decoding for CS-VAE vs. PS-VAE.

522 We introduced a model based on transfer learning to perform the decoding test on the previously tested
 523 subject. The rest of the three mice were the input to the original training model. The procedures were
 524 similar to before, after the trials were shuffled and split, we decoded the behavior directly with the raw neural
 525 activities with the time window being 0.15s. After that, we implemented three perceptron layers for each of
 526 the three mice before the output of which went into a recurrent neural network (RNN). The RNN consisted
 527 of one long short-term memory (LSTM) layer with a unit number of 64 and a drop-out layer with a rate
 528 being 0.5. We applied the Adam optimizer with learning rate decay from 0.1 with 0.3 decay rate for every 5
 529 step. The batch size was 150 and we trained for 200 epochs. After we finished training the original network,
 530 we transferred the RNN model to the new model which was applied to train the fourth mouse alone. For
 531 the fourth mouse, the trials were split with different training and testing ratios. After applying the same
 532 steps to the data, the neural activities then went through a new perceptron layer before going through the
 533 pre-trained RNN model. We applied the Adam optimizer with the same learning rate decay procedures as
 534 well. We again, trained for 200 epochs with batch size being 128 this time. The trade-off between accuracy
 535 and time for different models can be found in Tables 4 and 5.

Table 4: Training size vs R^2 value for multi-subject dataset

Training size	Linear model	Dense model	LSTM model	Transfer learning model
67712	0.476 ± 0.048	0.580 ± 0.058	0.610 ± 0.054	0.590 ± 0.050
58512	0.478 ± 0.014	0.560 ± 0.023	0.595 ± 0.022	0.579 ± 0.019
49312	0.483 ± 0.009	0.556 ± 0.014	0.593 ± 0.013	0.576 ± 0.011
40112	0.476 ± 0.013	0.543 ± 0.019	0.576 ± 0.019	0.562 ± 0.015
30912	0.470 ± 0.0011	0.529 ± 0.015	0.559 ± 0.018	0.552 ± 0.013
21712	0.458 ± 0.010	0.496 ± 0.016	0.524 ± 0.017	0.522 ± 0.013
12512	0.424 ± 0.012	0.461 ± 0.0019	0.480 ± 0.025	0.485 ± 0.018
3312	0.269 ± 0.030	0.321 ± 0.048	0.325 ± 0.057	0.345 ± 0.043

Table 5: Training size vs time usage for multi-subject dataset

Training size	Linear model	Dense model	LSTM model	Transfer learning model
67712	1.442 ± 0.282	115.946 ± 1.559	169.801 ± 5.961	169.482 ± 5.041
58512	1.130 ± 0.212	80.428 ± 1.586	146.734 ± 5.063	151.771 ± 4.162
49312	0.937 ± 0.194	68.879 ± 1.257	122.500 ± 2.283	125.240 ± 4.479
40112	0.679 ± 0.114	56.449 ± 1.119	100.336 ± 2.212	102.923 ± 3.391
30912	0.484 ± 0.079	44.427 ± 0.808	78.850 ± 1.771	79.907 ± 2.608
21712	0.309 ± 0.050	31.968 ± 0.574	57.670 ± 1.681	56.032 ± 1.549
12512	0.162 ± 0.008	19.573 ± 0.369	35.557 ± 1.050	33.409 ± 0.741
3312	0.104 ± 0.034	7.292 ± 0.092	13.318 ± 0.342	11.365 ± 0.336

.11 Multidimensional Canonical Correlation Analysis (MCCA) for neural signal alignment

In our work, after extracting similar behaviors chunks from different individuals, we then extracted the corresponding neural activity for each subject. To smooth away the discreteness of the neural activity chunks, we shuffled the chunks before concatenating them together. After that, we performed the MCCA for all four subjects on each brain region. For each brain region, we choose the four sets of neural activities being the same length d , $X1 = \{x1_1, x1_2, \dots, x1_n\} \in R^{n \times d}$, $X2 = \{x2_1, x2_2, \dots, x2_n\} \in R^{m \times d}$, $X3 = \{x3_1, x3_2, \dots, x3_n\} \in R^{k \times d}$, and $X4 = \{x4_1, x4_2, \dots, x4_n\} \in R^{l \times d}$. Here, we choose the minimum number of region dimensionality in all of the four subjects as the dimension of canonical coordinate space, $minimum\{n, m, k, l\}$, and is annotated as j . For each dimension, define the projection weights for each dataset as $a_j = \{a_{j1}, a_{j2}, \dots, a_{jn}\}$, $b_j = \{b_{j1}, b_{j2}, \dots, b_{jn}\}$, $c_j = \{c_{j1}, c_{j2}, \dots, c_{jn}\}$, and $d_j = \{d_{j1}, d_{j2}, \dots, d_{jn}\}$. The resulting projected datasets are now d -dimensional arrays: $u1_j = \langle a_j, X1 \rangle$, $u2_j = \langle b_j, X2 \rangle$, $u3_j = \langle c_j, X3 \rangle$, and $u4_j = \langle d_j, X4 \rangle$. For each of the coordinate spaces, the objective functions can be written as:

$$\rho_j = \frac{\langle u1_j, u2_j, u3_j, u4_j \rangle}{\|u1_j\| \|u2_j\| \|u3_j\| \|u4_j\|} \quad (13)$$

Generally, for each pair of canonical components, the above equation is solved iteratively to find the best projects that can maximize the correlation. During training, the orthogonality between each canonical component is constrained. In our experiment, we calculated the across-subject correlations for each obtained CCs and kept the highest correlation value for each pair, here termed ρ_1 (Equation 13). We performed the above task for each brain region. In addition, we shuffled the chunks ten times and repeated the above steps.

We also calculated the canonical component for the same subject having similar behaviors. We applied the same methods as stated above to find similar behavior components and the corresponding neural activities. We divided the obtained neural activities into two parts with the same length and performed the CCA on those two signals. We calculated the correlation between the first two canonical correlation axes as the baseline.

561 **.12 Code**

562 The code for training the CS-VAE can be found in Supplementary Material 7. The code can be executed
563 by simply compiling the script 'train.py'. All the code are available at: <https://github.com/saxenalab->
564 [neuro/Behavioral-feature-extraction-CS-VAE](https://github.com/saxenalab-neuro/Behavioral-feature-extraction-CS-VAE).