

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 in a unified manner 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 and the efficient decoding of a novel subject's behavior.

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, an 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 include directly applying dimensionality reduction methods such as Principal Component Analysis (PCA) and Variational Autoencoders (VAEs) to video data [2, 8, 9]. However, understanding or segmentation of the latent variables is difficult for any downstream tasks such as motif generation. To

combine the best of both worlds, semi-supervised VAEs have been used for the joint estimation of tracked body parts and unsupervised latents that can effectively describe the entire image [2]. These have not been applied to across-subject data, with the exception of [10], where the authors directly use a frame of each subject’s video as a context frame to define individual differences; however, this method only works with a *discrete* set of *labeled* sessions or subjects. These methods fail when applied without labeled subject data, or more importantly, when analyzing freely-behaving social behavior, due to continuously shifting image distributions that confound the latent space.

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

To further demonstrate the utility of our approach, we show that we can recover behavioral motifs from the resulting latents in a seamless manner. We recover (a) the same motifs across different animals performing the same task, and (b) motifs pertaining to social interactions in a freely moving task with two animals. Furthermore, we show the neural decoding of multiple animals in a unified model, with benefits towards the efficient decoding of the behavior of a novel subject.

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

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

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

2 Results

2.1 CS-VAE Model Structure

Although existing pose estimation methods are capable enough to capture the body position of the animals in both open and contained space, tracking specific actions such as shaking and wriggling still remains a problem. However, a purely unsupervised or semi-supervised model such as a VAE or PS-VAE lacks the

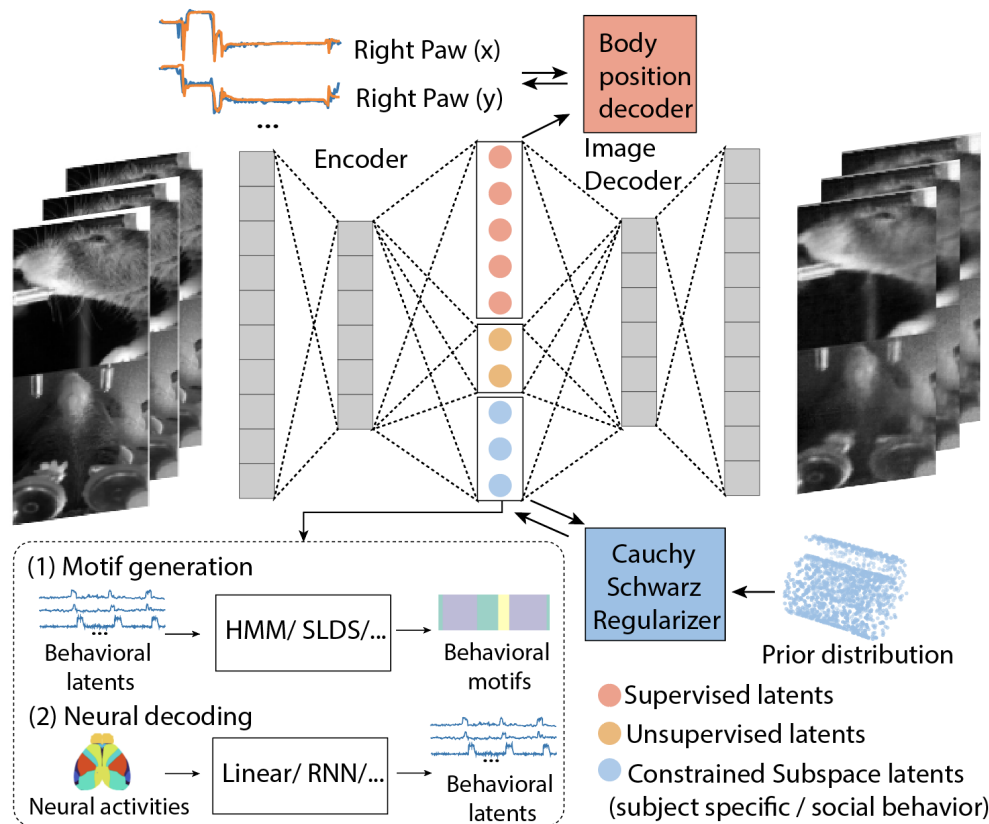


Figure 1: Overview of the Constrained Subspace Variational Autoencoder (CS-VAE). The latent space is divided in three parts: (1) the supervised latents decode the labeled body positions, (2) the unsupervised latents model the individual’s behavior that is not explained by the supervised latents, and (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.

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 to
 78 distinct parts of the subspace as in [2, 10], we directly constrain the latent subspace to a flexible prior dis-
 79 tribution 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
 85 supervised tracking methods such as DLC. The unsupervised latents capture the rest of the individual’s be-
 86 havior 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).

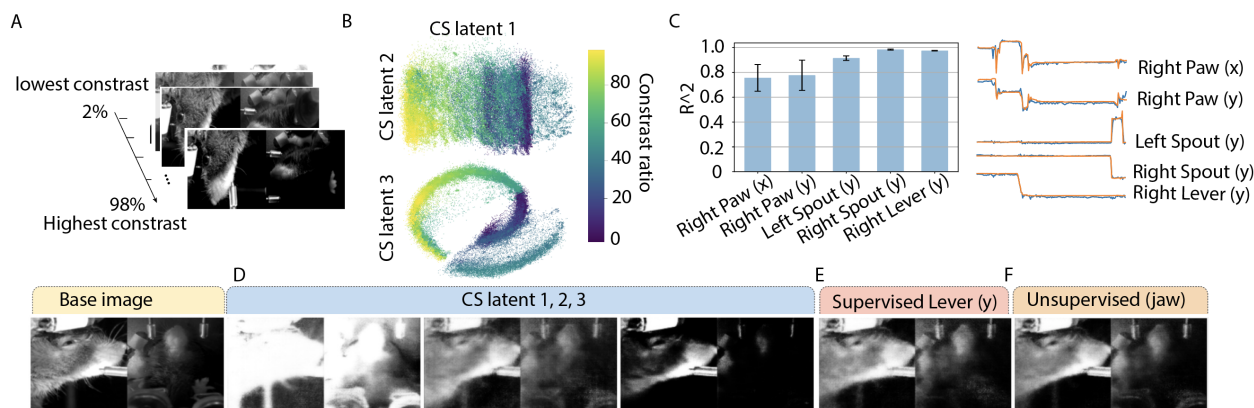


Figure 2: (A) Simulated dataset: behavioral videos from one mouse with artificially simulated differences in contrast. (B) Distribution occupied by the 3 CS latents. The constrained latents are distributed according to the pre-defined prior: a Swiss roll distribution. Different contrast ratios separate 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.

2.2 Modeling Smooth Variations in a Simulated Dataset

We performed a simulation study on the behavioral videos of one of the mice in the ‘Multi-Subject Behavior’ dataset detailed in Appendix A. We applied a continuously varying contrast ratio throughout the trials (Fig. 2A) to model smoothly varying lighting differences across the dataset. We then randomly shuffled all the trials and trained a CS-VAE model with a swiss roll as a prior distribution. Here, the R^2 for the supervised labels was 0.881 ± 0.05 (Fig. 2C), and the mean squared error (MSE) for reconstructing the entire frame was 0.0067 ± 0.0003 , showing that both the images and the labels were fit well. This was comparable to the PS-VAE model, where the R^2 for the supervised labels was 0.881 ± 0.09 , and the MSE for the entire frame reconstruction was $3.24 \pm 3.51 \cdot 10^{-5}$.

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

2.3 Modeling Multi-Subject Behavior

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

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

Disentangled Latent Space Representation We show latent traversals for each mouse in Fig. 4, with the base image chosen separately for each mouse (videos in Supplementary Material 3). We see that, even for different mice, the supervised latent can successfully capture the corresponding labeled body part

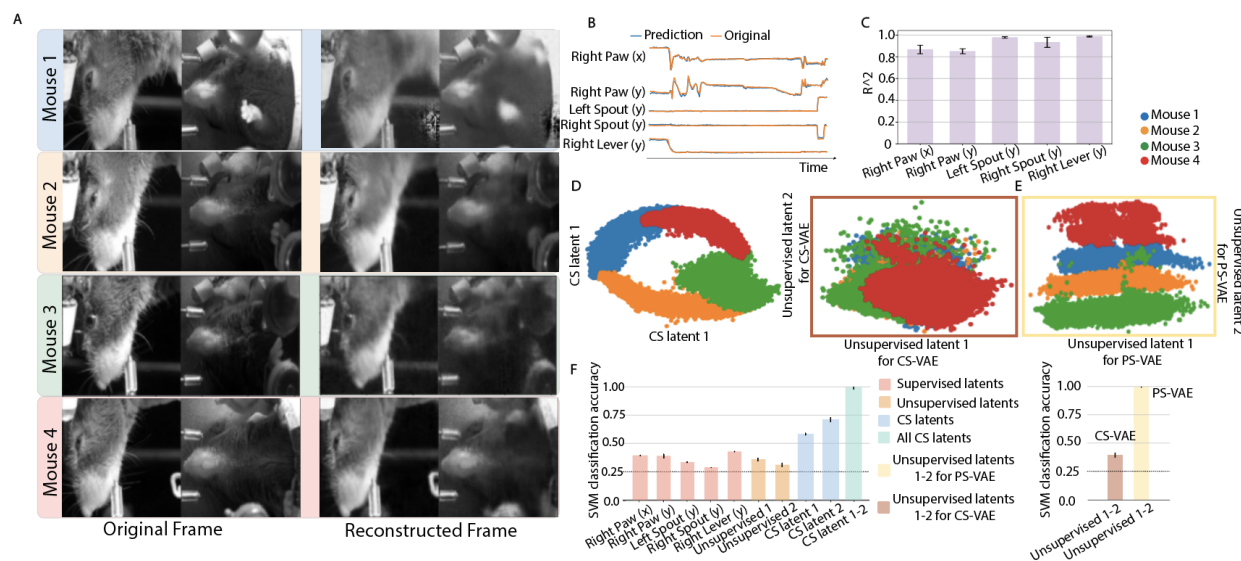


Figure 3: Modeling the behavior of four different mice. A. Image reconstruction result for an example frame from each mouse. B. Label reconstruction result for an example trial. C. R^2 value for label reconstruction for all mice. D. (Left) CS latent and (Right) unsupervised latent distributions for all 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.

(Fig. 4A). The example unsupervised latent is shown to capture parts of the jaw of each mouse (Fig. 4B), and is well-localized, comparable with the example supervised latent. The CS latent dimension encodes many different parts of the image, and has a large effect on the appearance of the mouse, effectively changing the appearance from one mouse to another, signifying that it is useful in the case of modeling mouse-specific differences (Fig. 4C). We demonstrate the abilities of the CS latent in capturing the appearance of the mouse by directly changing the CS latent from one part of subspace to another (Figure 4D). The changes in appearance along with the invariance in actions shows the intraoperability between mice by only changing the CS latents in this model (Fig. 4D).

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

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

Lastly, we trained the model while using prior distributions of different types, to understand the effect

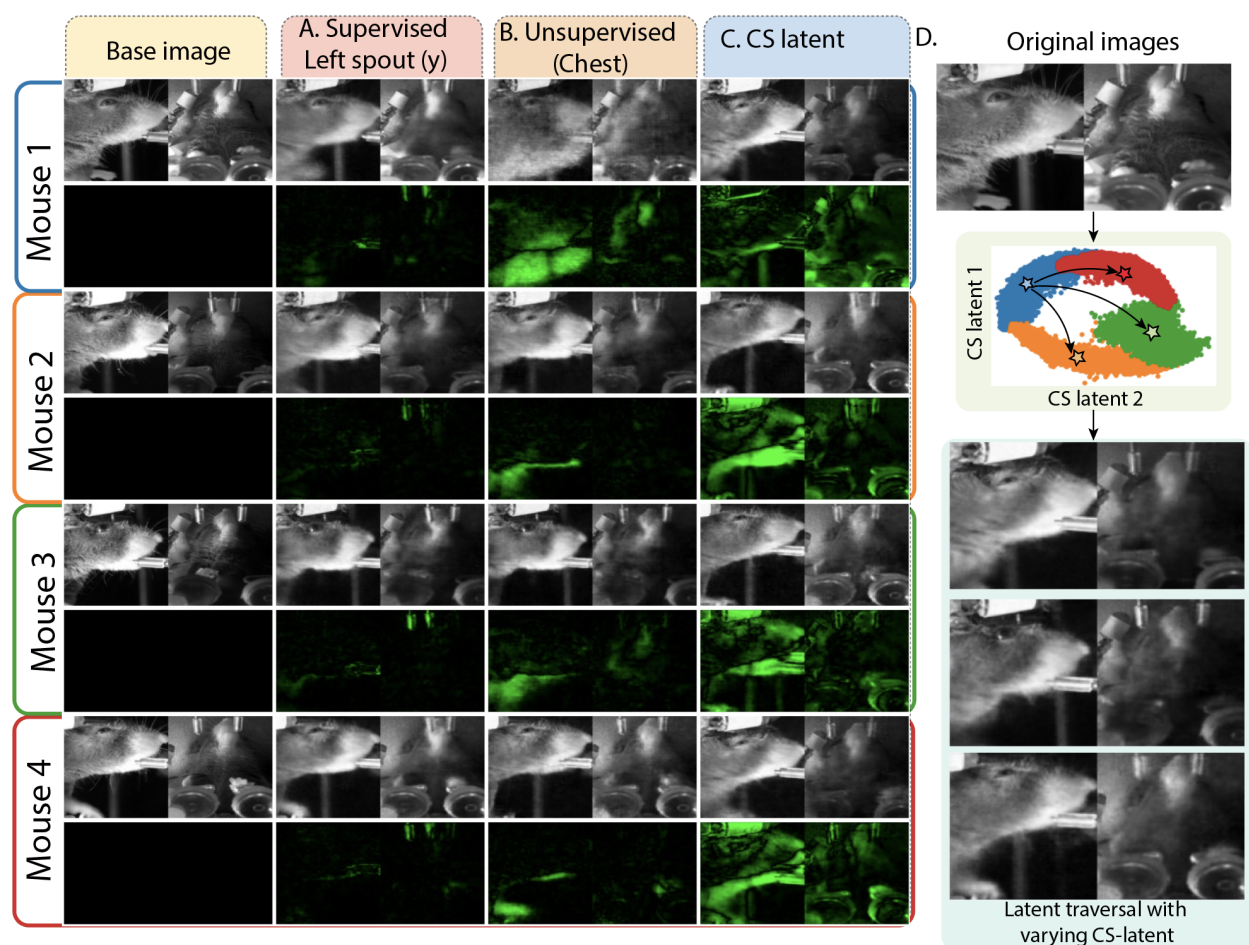


Figure 4: Latent traversals for behavioral modeling of four different mice for A. an example supervised latent that captures the left spout across all the subjects, B. an example unsupervised latent that captures the chest of the mice, and C. an example CS latent that successfully captures the mouse 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.

145 on the separability of the resulting latents. The separability was comparable across a number of different
 146 prior distributions, such as a swiss roll and a plane, signifying that the exact type of prior distribution
 147 does not play a large role.

148 **Across-Subject Motif Generation** To further show that the supervised and unsupervised latents
 149 produced by CS-VAE are interoperable between the different mice, we apply a standard SLDS model to
 150 uncover the motifs using this across-subject subspace. As seen in the ethograms (left) and the histograms
 151 (right) in Fig. 5, the SLDS using the CS-VAE latents captures common states across different subjects,
 152 indicating that the latents are well overlapped across mice. The supervised latents related to equipment
 153 in the experiment, here the spout and lever, split the videos into four states (different colors in the
 154 ethograms in Fig. 5A), that we could independently match with ground truth obtained from sensors in
 155 these equipment. The histograms show that, as expected, these states occur with a very similar frequency
 156 across mice. We also explored the behavioral states related to the right paw. The resulting three states
 157 captured the idle vs. slightly moving vs. dramatically moving paw (Fig. 5B). The histograms show that
 158 these states also occur with a very similar frequency across mice. Videos for all these states are available
 159 in Supplementary Material 2. Lastly, we extracted the behavioral states related to the unsupervised
 160 latents, which yielded 3 states related to raising of the paws (including grooming) and jaw movements
 161 (including licking) that are present in all four mice, as shown in Fig. 5C. We see that different mice have
 162 different tendencies to lick and groom, e.g., mouse 1 and 4 seem to groom more often.

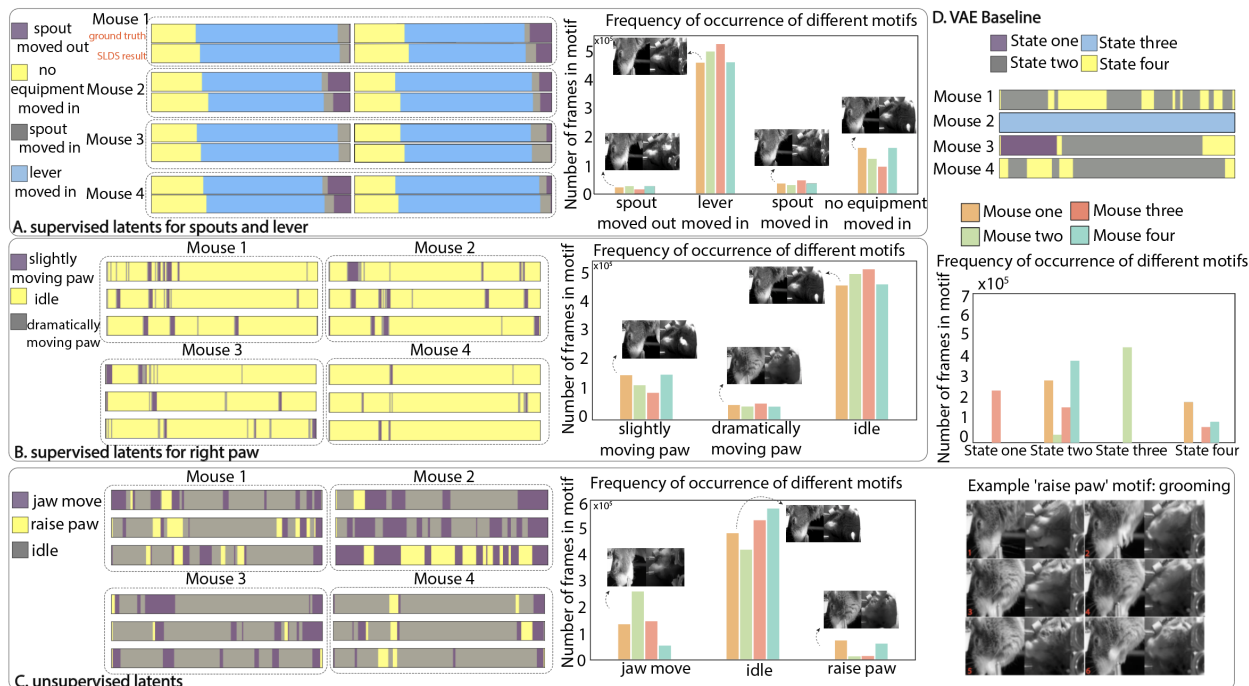


Figure 5: Motif generation for across-subject (supervised and unsupervised) behavioral latents using CS-VAE. SLDS results for CS-VAE latents: A. Supervised latents relating to equipment in the field of view. The equipment actions are similar for each trial. B. Supervised latents relating to tracked body 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 across-subject motifs.

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

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

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

2.4 Modeling Freely-Moving Social Behavior

The dataset consists of a 16 minute video of two adult novel C57BL/6J mice, a female and a male, interacting in a clean cage. Prior to the recording session the mice were briefly socially isolated for 15 minutes to increase interaction time. As preprocessing, we aligned the frame to one mouse and cropped the video (schematic in Fig. 7A; details in the Appendix B). We tracked the nose position (x and y coordinates) of the mouse using DLC. Here, we did not include an unsupervised latent space, since the alignment and supervised labels resulted in the entire individual being explained well using the supervised latents.

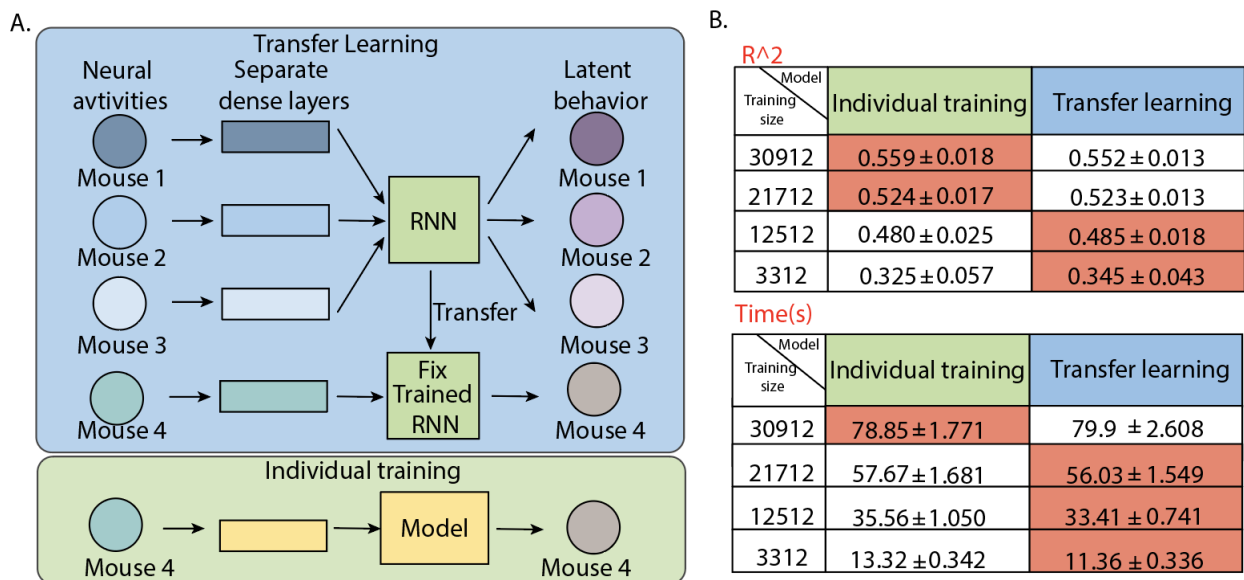


Figure 6: A. Transfer learning model framework. Each of the four mice has a specific dense layer for aligning the neural activities. After the model is trained using three mice, the across-subject Recurrent Neural Network (RNN) layer is fixed and transferred to the fourth mouse. As a comparison, we trained 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.

185 **Reconstruction Accuracy** The CS-VAE model results in a mean label reconstruction accuracy
 186 0.961 ± 0.0017 (Fig. 7B), with the MSE for frame reconstruction as $1.21 \cdot 10^{-5}$ (Fig. 7B). We compared
 187 the performance of our model with the VAE and PS-VAE (Table 1), and the CS-VAE model performed
 188 better than the baseline models for both image and label reconstruction. For the VAE, we obtained the
 189 R^2 for nose position prediction by training a multi-layer perceptron (MLP) with a single hidden layer
 190 from the VAE latents to the nose position.

191 **Disentangled Latent Space Representation** We calculated the latent traversals for each latent
 192 as in Section 4. As shown in the videos in Supplementary Material 4, CS latent 1 captures the second
 193 mouse to the front of the tracked mouse, CS latent 2 captures the front and above position of the second
 194 mouse, and CS latent 3 captures the position where the second mouse is below the tracked mouse.

195 To visualize the latent space and understand the relationship to social interactions, we plot the CS la-
 196 tents overlaid with the nose-to-tail distance between the two mice (nose of one mouse to the tail of the other)
 197 in Fig. 7C. We see that the CS latents represent the degree of social interaction very well, with a large separa-
 198 tion between different social distances. Furthermore, we trained an MLP with a single hidden layer from dif-
 199 ferent models' latents to the nose-to-tail distance, and the CS-VAE produces the highest accuracy (Table 1).

200 **Motif Generation** We applied a hidden Markov model (HMM) to the CS latents to uncover behav-
 201 ioral motifs. The three clusters cleanly divide the behaviors into social investigation vs. non-social
 202 behavior vs. non-social behavior with the aligned mice exploring the environment. To effectively visualize
 203 the changes in states, we show the ethogram in Fig. 8A. Videos related to these behavioral motifs are

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

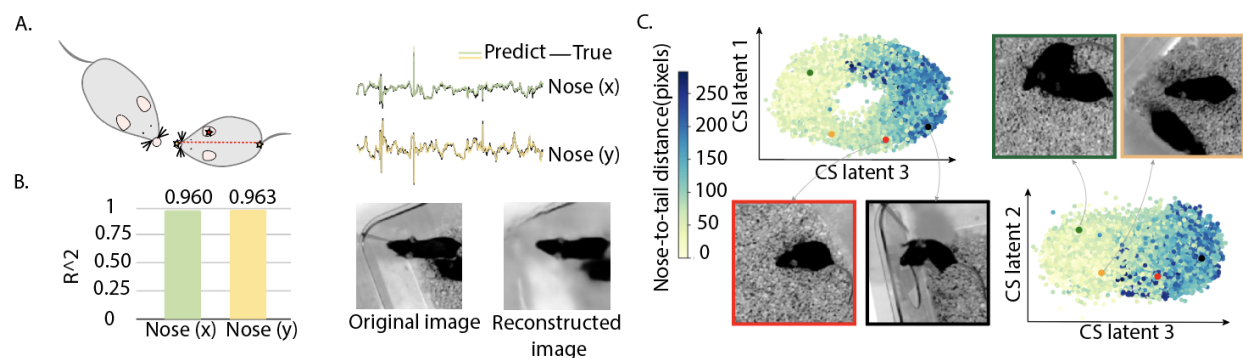


Figure 7: A. Image alignment for the social behavior data. B. Model performance on the social behavior dataset. C. Visualization of the CS latents overlaid with the nose-to-tail distance between the two interacting mice. The CS latents separates the frames that contain social interactions from those that do not.

204 provided in Supplementary Material 5.

205 Lastly, we calculated different metrics to quantitatively evaluate the difference between each behavioral
 206 motif. The results are shown in Fig 8B, where we plot the average values for distances and angles between
 207 different key points. The lower distance between the two mice in *State a* demonstrates that the mice
 208 are close to each other in that state, pointing to social interactions. The smaller nose-to-tail distance
 209 for the aligned mouse in *State c* points to this state encoding for the ‘rearing’ of the mouse. The angle
 210 between the two mice further reveals the relative position between the two mice; in *State b*, the second
 211 mouse is located above the aligned mouse, while the opposite is true for *State c*. These metrics uncover
 212 the explicit differences between the different motifs that are discovered by CS-VAE.

213 3 Discussion

214 In the field of behavior modeling, there exist three major groups of methods, supervised, unsupervised,
 215 and semi-supervised. The supervised methods consist of methods such as DeepLabCut (DLC) [7], LEAP
 216 [6], AlphaTracker [5], amongst others. Although these methods capture the positions of the subjects, they
 217 lack the ability to model smaller movements and unlabeled behavior, and necessitate tedious labeling. On
 218 the other hand, unsupervised methods such as MoSeq [9] and Behavenet [8] lack the ability to produce
 219 interpretable behavioral latents. While some semi-supervised methods, for instances, MSPS-VAE [2]
 220 and DBE [10], succeed in producing interpretable latents and modeling behavior across subjects, they
 221 need significant human input, and lack the ability to model freely-moving animals’ behavior. Here, we
 222 introduce a constrained generative network called CS-VAE that effectively addresses major challenges
 223 in behavioral modeling- disentangling multiple subjects and representing social behaviors.

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

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

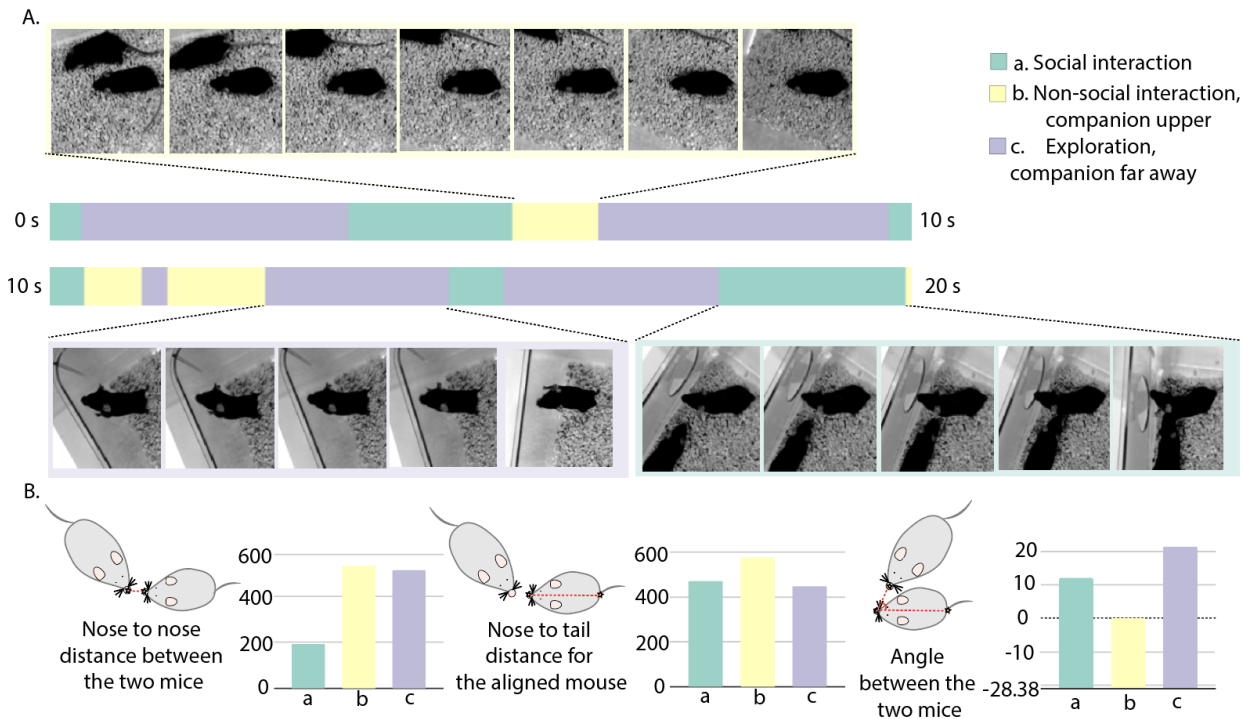


Figure 8: A. Ethogram for the animals’ behavior recovered using hidden Markov models (HMM) applied to the CS latents. B. Different metrics for analysing the behavioral motifs. Here, the three motifs are *a*. social interaction; *b*. non-social interaction with the companion on the upper side of the aligned 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.

4 Methods

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

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

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

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

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

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 the estimated quadratic entropy of $p_1(x)$ and $p_2(x)$, respectively, while $-\log(V(p_1, p_2))$ is the estimated 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 Gaussian. N , N_1 , and N_2 are the number of samples being input into the model while σ is the kernel size. The choice of the kernel size depends on the dataset itself; generally, the kernel size should be greater than the number of the groups in the data. Equation (1) can be expressed as:

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

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

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

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

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

261 Here, the terms \mathcal{L}_{frames} and \mathcal{L}_{label} represent the reconstruction loss of the frames and the labels, respec-
262 tively. The \mathcal{L}_{KL-s} represents for the KL-divergence loss for the supervised latents while \mathcal{L}_{ICMI} , \mathcal{L}_{TC} ,
263 and \mathcal{L}_{DWKL} form the decomposed version of the KL loss for the unsupervised latents. Lastly, the \mathcal{L}_{CS}
264 represents the CS-divergence loss on our constrained latents. α is introduced to control the reconstruction
265 quality of the labels, β is adopted to assist the model in producing independent unsupervised latents,
266 and γ is implemented to control the variability in the constrained latent space for better separation.
267 The detailed explanations and derivations for each term in the objective function are in Appendix C.
268 Furthermore, the loss terms in Equation (5) can be modified to fit various conditions. For a freely-behaving
269 social task, the background for one individual in the container could be the edge of the container as well
270 as the rest of the individuals in the container. The choice of hyperparameters and the loss curves through
271 the training process is shown in Appendix E and G, respectively.

272 **Visualization of the latent space** To test how the image varies with a change in the latent, one
273 frame from the trials is randomly chosen as the ‘base image’, and the effect of varying a specific latent
274 at a time is visualized and quantified. This is known as the ‘latent traversal’ [2]. First, for each latent
275 variable, we find out the maximum value that it occupies across a set of randomly selected trials. We
276 then change that specific latent to achieve its maximum value, and this new set of latents forms the input
277 to the decoder. We obtain the corresponding output from the decoder as the ‘latent traversal’ image.
278 Finally, we visualize the difference between the ‘latent traversal’ image and the base image. The above
279 steps are performed for each latent individually. In videos containing latent traversals (Supplementary
280 Material), we change the latent’s value from its minimum to its maximum across all trials, and input
281 all the corresponding set of latents into the decoder to produce a video.

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

289 **Efficient Neural Decoding** Decoding neural activity to predict behavior is very useful in the un-
290 derstanding of brain-behavior relationships, as well as in brain-machine interface tasks. However, models
291 to predict high-dimensional behavior using large-scale neural activity can be computationally expensive,
292 and require a large amount of data to fit. In a task with multiple subjects, we can utilize the similarities
293 in brain-behavior relationships to efficiently train models on novel subjects using concepts in transfer
294 learning. Here, we represent across-subject behavior in a unified manner and train an across-subject neural
295 decoder. Armed with this across-subject decoder, we show the decoding power on a novel subject with
296 varying amounts of available data, such that it can be used in a low-data regime. The implementational
297 details for this transfer learning approach can be found in Appendix J.

298 References

- 299 [1] Krakauer, J. W., Ghazanfar, A. A., Gomez-Marin, A., MacIver, M. A. Poeppel, D. Neuroscience
300 needs behavior: Correcting a reductionist bias. *Neuron* 93, 480–490 (2017).
301 [2] Whiteway, M. R. et al. Partitioning variability in animal behavioral videos using semi-supervised
302 variational autoencoders. *bioRxiv* (2021).
303 [3] Mathis, A. et al. Deeplabcut: markerless pose estimation of user-defined body parts with deep
304 learning. *Nature Neuroscience* 21, 1281–1289 (2018).

- 305 [4] Pereira, T. et al. Fast animal pose estimation using deep neural networks. bioRxiv (2018).
306 [5] Chen, Z. et al. Alphatracker: A multi-animal tracking and behavioral analysis tool. bioRxiv (2020).
307 [6] Pereira, T. D. et al. Publisher correction: Slep: A deep learning system for multi-animal pose
308 tracking. Nat Methods (2022).
309 [7] Lauer, J. et al. Multi-animal pose estimation and tracking with deeplabcut. bioRxiv (2021).
310 [8] Batty, E. et al. Behavenet: nonlinear embedding and bayesian neural decoding of behavioral videos.
311 In Wallach, H. et al. (eds.) Advances in Neural Information Processing Systems, vol. 32 (Curran
312 Associates, Inc., 2019).
313 [9] Wiltschko, A. B. et al. Revealing the structure of pharmacobehavioral space through motion
314 sequencing. Nature neuroscience 23, 1433–1443 (2020).
315 [10] Shi, C. et al. Learning disentangled behavior embeddings. In NeurIPS (2021).
316 [11] Santana, E., Emigh, M. Principe, J. Information theoretic-learning auto-encoder (2016).
317 [12] Tran, L., Pantic, M. Deisenroth, M. P. Cauchy-schwarz regularized autoencoder (2021). 2101.02149.
318 [13] Wiltschko, A. et al. Mapping sub-second structure in mouse behavior. Neuron 88, 1121–1135 (2015).
319 [14] Musall, S., Kaufman, M. T., Juavinett, A. L., Gluf, S. Churchland, A. K. Single-trial neural
320 dynamics are dominated by richly varied movements. Nature neuroscience 22, 1677–1686 (2019).
321 [15] Saxena, S. et al. Localized semi-nonnegative matrix factorization (locanmf) of widefield calcium
322 imaging data. PLOS Computational Biology 16, 1–28 (2020).
323 [16] Linderman, S. et al. Bayesian Learning and Inference in Re- current Switching Linear Dynamical
324 Systems. In Singh, A. Zhu, J. (eds.) Proceedings of the 20th International Conference on Artificial In-
325 telligence and Statistics, vol. 54 of Proceedings of Machine Learning Research, 914–922 (PMLR, 2017).