



OPEN ACCESS

EDITED BY

Shyam Diwakar,
Amrita Vishwa Vidyapeetham, India

REVIEWED BY

Adam Ponzi,
National Research Council, Italy
Remi Monasson,
École Normale Supérieure, France

*CORRESPONDENCE

Haoming Yang
✉ haoming.yang@duke.edu

[†]These authors have contributed equally to this work

RECEIVED 22 September 2025

REVISED 27 January 2026

ACCEPTED 28 January 2026

PUBLISHED 13 February 2026

CITATION

Yang H, Angjelichinoski M, Wu S, Putney J, Sponberg S and Tarokh V (2026) Cross-subject mapping of neural activity with restricted Boltzmann machines.

Front. Comput. Neurosci. 20:1710914.
doi: 10.3389/fncom.2026.1710914

COPYRIGHT

© 2026 Yang, Angjelichinoski, Wu, Putney, Sponberg and Tarokh. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Cross-subject mapping of neural activity with restricted Boltzmann machines

Haoming Yang^{1*†}, Marko Angjelichinoski^{1†}, Suyu Wu¹,
Joy Putney², Simon Sponberg² and Vahid Tarokh¹

¹Department of Electrical and Computer Engineering, Duke University, Durham, NC, United States,

²School of Physics and Biological Sciences, Georgia Institute of Technology, Atlanta, GA, United States

Subject-to-subject variability is a common challenge in generalizing neural data models across subjects, discriminating subject-specific and inter-subject features in large neural datasets, and engineering neural interfaces with subject-specific tuning. While many methods exist that map one subject to another, it remains challenging to combine many subjects in a computationally efficient manner, especially with highly non-linear features such as populations of spiking neurons or motor units. Consider subjects with trained neural decoders as sources and those without as targets. Our objective is to transfer data from one or more target subjects to the domain of the source subjects to directly apply the source neural decoder such that no target decoder needs to be trained. We propose to use the Restricted Boltzmann Machine (RBM) with Gaussian inputs and Bernoulli hidden units; once trained over the entire feature set of subjects, the RBM allows the mapping of target features on source feature spaces using Gibbs sampling. We also consider a novel computationally efficient training technique for RBMs based on the Fisher divergence, which allows closed-form gradients of the RBM to be computed. We apply our methods to decode turning behaviors from neuromuscular recordings of spike trains from the ten muscles that primarily control wing motion in an agile flying hawk moth, *Manduca sexta*. The dataset consists of this comprehensive motor program recorded from nine subjects, each driven by six discrete visual stimuli. The evaluations show that the target features can be decoded using the source classifier to classify the visual stimuli with an accuracy of up to 95% when mapped using an RBM trained by Fisher divergence, suggesting that RBMs for multi-cross-subject mapping applications are effective and efficient.

KEYWORDS

cross-subject mapping, distribution alignment, domain adaptation, restricted Boltzmann machine, transfer learning

1 Introduction

Combining data across individuals or subjects is a ubiquitous need in analyzing neuroscience experiments. Individual variation can produce subject-specific variation or noise and different subjects may utilize different strategies to represent the same sensory input or execute the same motor task. Yet, implementing reliable *cross-subject* algorithms in neuroscience is a notoriously challenging problem. Another important factor contributing to its difficulty arises from the *non-stationary* nature of the neural activity signals, whose statistical

properties vary dramatically even under slight changes in the recording conditions (Rao, 2013; Jayaram et al., 2016; Dabagia et al., 2023). As a result, the algorithms trained and optimized on data collected from a given subject, have long failed to perform reliably when directly applied to other subjects (Torres-Oviedo and Ting, 2010). Even modern neural decoders trained on one subject will perform close to a random choice classifier if applied directly to a different subject, thus failing to identify the correct neurological state or stimulus condition even when the subjects perform the same tasks simultaneously (Angelichinoski et al., 2020a).

State-of-the-art cross-subject mapping methods in neuroscience are commonly tackled through distribution alignment (Dyer et al., 2017; Lee et al., 2019). These approaches aim to find a mapping that aligns the distribution of target features to the distribution of source features. They generally require strong assumptions, which ensures their provably robust, accurate result. For example, Procrustes alignment-based methods require both source and target to exist during test time and are a fundamentally linear transformation (Haxby et al., 2011; Degenhart et al., 2020). Another example, the Hierarchical Wasserstein Alignment (HiWA) assumes the alignment can be solved through a linear transformation (Lee et al., 2019), and requires cluster structure in the dataset with cluster labels to solve the optimal transport optimization. These alignment-based methods perform well in many decoding problems, such as aligning the data of the same subject across time and distributional shifts (Karpowicz et al., 2022; Dyer et al., 2017; Lee et al., 2019). However, in a more complex multi-subject setting where relationships across subjects could be highly non-linear, the strong assumptions of alignment methods restrict their flexibility. These alignment methods could also be inefficient with high-dimensional datasets, as the matrix operations become more computationally intensive as dimensionality grows. Finally, the alignment based methods also require the source data during test time, which could be private, especially in clinical or commercial scenarios (CTRL-labs at Reality Labs et al., 2024).

While such distribution alignment methods are already adopted in neuroscience research, there remains a need for reliable, non-linear cross-subject mapping methods that are more flexible against complex real world data. Problems of this type, i.e., problems where the training and test data originate from different distributions, are common in deep learning and are typically tackled within the sub-field of *transfer learning*. Transfer learning methods are already starting to have a significant impact on neuroscience and behavior: for example, a reliable cross-neuron population mapping in the latent domain of neural activities can stabilize and improve brain-computer interface (Pandarinath et al., 2018; Degenhart et al., 2020; Karpowicz et al., 2022). In the context of the cross-subject problem, various approaches have already been considered (Jayaram et al., 2016). One direction to map between subjects is domain adaptation methods that map the target data onto the source feature spaces (Harvey et al., 2024; Angelichinoski et al., 2020a). Generative modeling is another promising approach to cross-subject mapping: one can apply a directed graph such as a conditional variational autoencoder (cVAE) to generate the target data onto the feature space of the source data (Angelichinoski et al., 2020b), where the learning model for the downstream task (e.g. classification or regression) is trained. Such an approach,

however, requires a separate directed graphical model to be trained each time a new downstream task and/or a new source subject is considered. This renders the already demanding deep learning method inefficient, especially when the number of models can scale exponentially as the number of subjects grows.

In this paper, we propose an efficient, non-linear algorithm that uses undirected graphs to generate samples for cross-subject mapping. Specifically, we leverage the non-linear generative algorithm known as the Restricted Boltzmann Machine (RBM) (Hinton, 2002; Carreira-Perpinan and Hinton, 2005), which allows us to bring the flexible and assumption-less deep learning perspectives and the versatility of the alignment-based methods together to construct an efficient cross-subject mapping method independent of its downstream decoding tasks. RBM is a popular generative model that has had notable success in representation learning with applications in a wide variety of tasks in neuroscience, including neuroimaging, classification of temporal events, and EEG analysis (Plis et al., 2014; Li et al., 2015; Hajinoroozi et al., 2016; Chai et al., 2017; Kim et al., 2020). RBMs have also found success in transfer learning, but mainly in computer vision applications (Wei and Pal, 2011; Wu and Ji, 2016; Farahani et al., 2020). The applications of RBM in neuroscience and transfer learning motivate us to adapt RBM for cross-subject mapping, which can serve as a general method with few assumptions on the dataset.

In the neuroscience context, we tackle the general problem of finding a common latent representation of high dimensional spiking data from many neurons across many individual subjects and mapping new subjects to this latent domain. While many other methods have aligned one target to one subject, few have considered the more general problem of combining representations across many subjects even though this is common to many experiments (Hajinoroozi et al., 2016; Kim et al., 2020). Only a few of the current methods consider highly non-linear representations such as spikes but only apply them in a naive environment where only two subjects exist and source and targets both exist during test time (Lee et al., 2019; Degenhart et al., 2020). Here we utilize RBMs to learn distribution mappings either from one source subject applied to many test subjects or from a combination of many source subjects applied to a single target. Therefore, during testing, the neural decoder trained on source data can be directly applied to decode the test data of the target subjects because the cross-subject mapping process is independent of downstream tasks. This particular feature of the undirected graphs allows efficient transfer between all subjects with one model, in comparison of the multi-model method induced by traditional subject mapping.

One limitation of RBMs is that they often are computationally expensive because contrastive divergence minimization does not have a closed form and requires Gibbs sampling. Here, we also extend the conventional contrastive divergence training of RBM (which is equivalent to maximum likelihood, we refer more details to Hinton, 2002) by proposing an alternative training based on Fisher divergence minimization (Hyvärinen, 2005). The Fisher divergence minimization allows the gradient of the RBM to be computed in closed form, fostering an even more efficient implementation that does not require iterative Gibbs sampling during training.

We evaluate the performance of our method for cross-subject decoding of discrete visual stimulus conditions using the spiking activity of the motor program, specifically the set of spiking motor units, in nine hawk moths (Putney et al., 2021). Each moth is exposed to the same set of six visual stimuli and the neuromuscular activity is collected in the form of spike trains extracted from fine wire electromyography (EMGs) of the ten primary flight muscles that control the wings, resulting in a comprehensive, spike-resolved motor program (Putney et al., 2019). Unlike vertebrate EMGs, these flight muscles act as effectively single motor units and result in identifiable spike trains comparable to population recordings of individual units elsewhere in the brain or peripheral nervous system.

In the context of these specific data, the challenge is to enable the decoding of a new target subject moth's motor program into the motor outcome (turn direction) based on a set of latent space shared in common with the other subjects (i.e. individual moths). Within a single subject, we can already nearly perfectly decode behavior as expected from a comprehensive motor program recording. However, the latent structures within population recordings across the different moths are different, so naive decoding does not generalize. Using our RBM with Fisher divergence we show that we can learn a generalized latent domain from the population of spiking units. We compare the decoding accuracy of this method to other common transfer learning methods most notably HiWA and RBMs with more traditional contrastive divergence. Our results demonstrate the promising potential of the proposed framework, with respectively up to 90% and 95% accuracy in decoding the behavioral state (i.e., the visual stimulus), when using the RBM trained with classical and the new Fisher divergence-based methods; these results thoroughly improve from previous state-of-the-art distribution alignment-based methods like HiWA.

2 Material and methods

We divide this section into three parts. In Section 2.1 we present the statistical formulation of the problem of cross-subject mapping as a problem of learning joint distribution between target and source feature vectors. Next, in Section 2.2, we discuss RBMs and present both the contrastive divergence and Fisher divergence training methods. Section 2.3 presents a simple cross-subject mapping approach that uses Gibbs sampling to draw samples from the joint distribution of the target and source features.

The code to reproduce this work is available at <https://github.com/imkeithyang/RBM-Cross-Subject-Mapping>.

2.1 Problem statement

We consider a case where the motor intentions are only known for the source subject while the population recordings are available for all target and source subjects, which requires the usage of cross-subject mapping during test time. The objective of cross-subject learning is to map target subject features to the feature space of

source subjects so that the motor intentions of the target subject can be decoded by the neural decoder that is learned on source subjects. In other words, we aim to obtain the appropriate source space representation of the tasks that the target subjects perform. Technically, the problem boils down to finding a function that maps the feature vectors. This mapping function can be assumed to be purely deterministic. However, in this paper, we adopt a probabilistic approach, which generates task-specific features in source feature space. We outline details below.

Let \mathcal{M}_T and \mathcal{M}_S denote the index sets of the target and source subjects, respectively. Aligning with standard transfer learning terminologies, source subject(s) is defined as the subject(s) where the neural decoder is available during test time, and target subject(s) is the subject that we wish to decode using the source decoder. Further, we let \mathcal{M} denote the joint set of all subjects, namely $\mathcal{M} = \mathcal{M}_T \cup \mathcal{M}_S$. For simplicity, we assume that the subsets \mathcal{M}_T and \mathcal{M}_S are disjoint, namely $\mathcal{M}_T \cap \mathcal{M}_S = \emptyset$. Let $\mathbf{x}_m \in \mathbb{R}^{D_m}$ denote the D_m -dimensional vector representing the neural activity of subject $m \in \mathcal{M}$; we refer to \mathbf{x}_m as the *feature vector* of subject m . Furthermore, we use $\mathbf{x}_T = (\mathbf{x}_i)_{i \in \mathcal{M}_T}$ and $\mathbf{x}_S = (\mathbf{x}_j)_{j \in \mathcal{M}_S}$ to respectively denote the joint feature vectors of the target and source subjects. Note that \mathbf{x}_T and \mathbf{x}_S are vectors with dimensions $D_T = \sum_{i \in \mathcal{M}_T} D_i$ and $D_S = \sum_{i \in \mathcal{M}_S} D_i$, respectively. Finally, we use \mathbf{x} to denote the joint vector of features of *all* subjects in the set \mathcal{M} ; that is $\mathbf{x} = (\mathbf{x}_T, \mathbf{x}_S) = (\mathbf{x}_m)_{m \in \mathcal{M}}$, and its dimension $D = \sum_{m \in \mathcal{M}} D_m$.

To learn the cross-subject mapping, we consider a *conditional probability distribution* to generate feature representations in the feature space of source subjects given the feature vector of target subjects. One option is to directly parameterize the probability density function $p(\mathbf{x}_S|\mathbf{x}_T)$ of this conditional distribution. Another approach is to first learn the distribution $p(\mathbf{x}_T)$ and the joint distribution $p(\mathbf{x}_S, \mathbf{x}_T)$ of all feature vectors across the entire population of subjects in \mathcal{M} , and then to obtain the conditional distribution $p(\mathbf{x}_S|\mathbf{x}_T)$ by Bayes' theorem.

For the purpose of cross-subject mapping, it is not mandatory to learn an explicit probability density function (pdf) of the conditional distribution. Recall that the objective is to obtain the feature representations of the target subject in the feature space of the source subject. Therefore, to this end, we aim to learn a generative framework such that we can easily sample from the conditional distributions $p(\mathbf{x}|\mathbf{h}; \theta)$ and $p(\mathbf{h}|\mathbf{x}; \theta)$ with \mathbf{h} denoting hidden variables. Generative algorithm for cross subject mapping in previous works, generally applies directed graphs such as cVAE (Sohn et al., 2015; Angelichinoski et al., 2020b), which map the target feature onto the feature space of the source subjects by the decoder of the autoencoder. However, directed graphs are not flexible to adopt new source subjects, and the architecture of cVAE is built on deep neural networks which are not easy to fine-tune given the limited size of data. In this paper, we consider the class of undirected graphical models and easily adapt to the generative framework. Given its relatively simple architecture and straightforward sampling scheme, an RBM is the first generative model we will tackle. In our application, the source data is unknown during test time, and our following construction of the RBM allows the mapping of multiple subjects from targets to sources with only one trained model.

This RBM is constructed such that a joint distribution of all features is learned. Specifically, we consider a generative model $p(\mathbf{x})$ of the concatenated feature vector $\mathbf{x} = (\mathbf{x}_T, \mathbf{x}_S)$. We first learn the joint distribution $p(\mathbf{x}) = p(\mathbf{x}_T, \mathbf{x}_S)$. Once the joint distribution is learned by the RBM, we use Gibbs sampling to sample from $p(\mathbf{x})$ with hidden variables \mathbf{h} . The Gibbs sampler proceeds as follows: we initialize the visible variables by $\mathbf{x}^{(0)} = (\mathbf{x}_T^{(0)}, \mathbf{x}_S^{(0)})$. Here, $\mathbf{x}_T^{(0)} = \mathbf{x}_T$ is given, and $\mathbf{x}_S^{(0)}$ is a dummy and noise-like vector. Next, we sample hidden variables $\hat{\mathbf{h}}$ from $p(\mathbf{h}|\mathbf{x}; \theta)$, and then sample visible variables $\mathbf{x}^{(1)}$ from $p(\mathbf{x}|\mathbf{h}; \theta)$. After sufficient sampling iterations, we expect to obtain $\mathbf{x}^{(k)} = (\mathbf{x}_T^{(k)}, \mathbf{x}_S^{(k)})$ where the $\mathbf{x}_S^{(k)}$ is the feature representation of \mathbf{x}_T in the feature space of source subjects. The hidden layer \mathbf{h} is designed to bridge target and source feature vectors. Note that this architecture setup allows us to only train one RBM that covers the need for any multi-subject mapping during test time. This flexibility follows from learning the joint distribution during training, allowing us to efficiently apply Bayes rule during sampling. During the training of the RBM, there is no distinction between \mathbf{x}_T and \mathbf{x}_S as all subjects are available, but the neural decoder can only be trained with \mathbf{x}_S as only the motor intention for the \mathcal{M}_S is available. Then during test time we decode the features of \mathcal{M}_T using the existing \mathcal{M}_S 's neural decoder. With reliable conditional sampling from the learned conditional distribution, we can now map new data from target to source without any source data available during testing, and apply the sampled data to the source's neural decoder. We illustrate this generative framework in Figure 1.

2.2 Learning restricted Boltzmann machines

We will first discuss the Gauss–Bernoulli RBMs and outline the principles of their training. We then review the standard training technique that aims to minimize the contrastive divergence, which is equivalent to minimizing the Kullback–Leibler (KL) divergence between the data-generating distribution and the model (Hinton, 2002; Carreira-Perpinan and Hinton, 2005). We also consider an alternative training technique that minimizes the Fisher divergence (Hyvärinen, 2005). Unlike the classical method, minimizing the Fisher divergence approach allows the gradients of the loss function to be computed in closed form. This improves the computational efficiency and reliability of the training.

2.2.1 Notation

We adopt the following notation conventions. Recall that in the context of the cross-subject mapping problem outlined in Section 2.1, the vector \mathbf{x} comprises the feature vectors of the entire population of subjects, i.e., $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_M)$. Let $\nabla_{\mathbf{x}}$ and $\Delta_{\mathbf{x}}$ denote the gradient and Laplacian operator with respect to (w.r.t.) the vector \mathbf{x} . Let $\text{dg}(\mathbf{x})$ denote a diagonal matrix whose main diagonal is \mathbf{x} . For a square matrix \mathbf{A} , let $\text{dg}(\mathbf{A})$ denote a diagonal matrix formed by setting all the elements to \mathbf{A} not on the main diagonal to zeroes. We use $\|\mathbf{x}\|$ (respectively $\|\mathbf{A}\|$) to denote the L_2 -norm of the \mathbf{x} (respectively the Frobenius norm of \mathbf{A}). We further use $p_*(\mathbf{x})$ to denote the true data-generating distribution of \mathbf{x} . In practice,

$p_*(\mathbf{x})$ is usually unknown; therefore, given a set of observations, a standard problem is to estimate the model density $p(\mathbf{x})$ from some model class that best explains the data under an appropriate evaluation metric. In this paper, we focus on a parametric density model class $p(\mathbf{x}; \theta), \theta \in \Theta$, which is parameterized as an RBM (see Section 2.2.4).

2.2.2 Kullback–Leibler divergence and the logarithmic loss

A common practice to measure the deviation of a postulated probability distribution $p(\mathbf{x})$ from the data-generating distribution $p_*(\mathbf{x})$ is to use the KL divergence defined by

$$D_{\text{KL}}(p_*, p) = -\mathbb{E}_* [\log p(\mathbf{x}; \theta)] + \mathbb{E}_* [\log p_*(\mathbf{x})], \quad (1)$$

where the expectation is taken w.r.t. $p_*(\mathbf{x})$ (as denoted by the subscript). $D_{\text{KL}}(p_*, p) \geq 0$ with equality if and only if $p = p_*$ almost surely. The minimization of $D_{\text{KL}}(p_*, p)$ is equivalent to the minimization $\mathbb{E}_* [\ell(\mathbf{x}; \theta)]$ where $\ell(\mathbf{x}; \theta) = -\log p(\mathbf{x}; \theta)$ is called the logarithmic loss.

Let θ_* denote the data-generating parameter that minimizes the KL divergence in Equation 1, namely $p(\mathbf{x}; \theta_*)$ is closest to $p_*(\mathbf{x})$ among all distributions over Θ under the KL divergence. It can be shown from the law of large numbers and standard regularity conditions that the maximum likelihood estimate (MLE) $\hat{\theta}_{\text{ML}} = \arg \min_{\theta} \bar{\ell}(\mathbf{x}; \theta)$ satisfies $\hat{\theta}_{\text{ML}} \rightarrow \theta_*$ in probability as the number of data points grows (White, 1982). In other words, the MLE is consistent.

2.2.3 Fisher divergence and Hyvärinen score

The Fisher divergence of $p(\mathbf{x}, \theta)$ from the data-generating pdf $p_*(\mathbf{x})$ is defined by

$$D_{\text{F}}(p_*, p) = \frac{1}{2} \mathbb{E}_* [\|\nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta) - \nabla_{\mathbf{x}} \log p_*(\mathbf{x})\|^2], \quad (2)$$

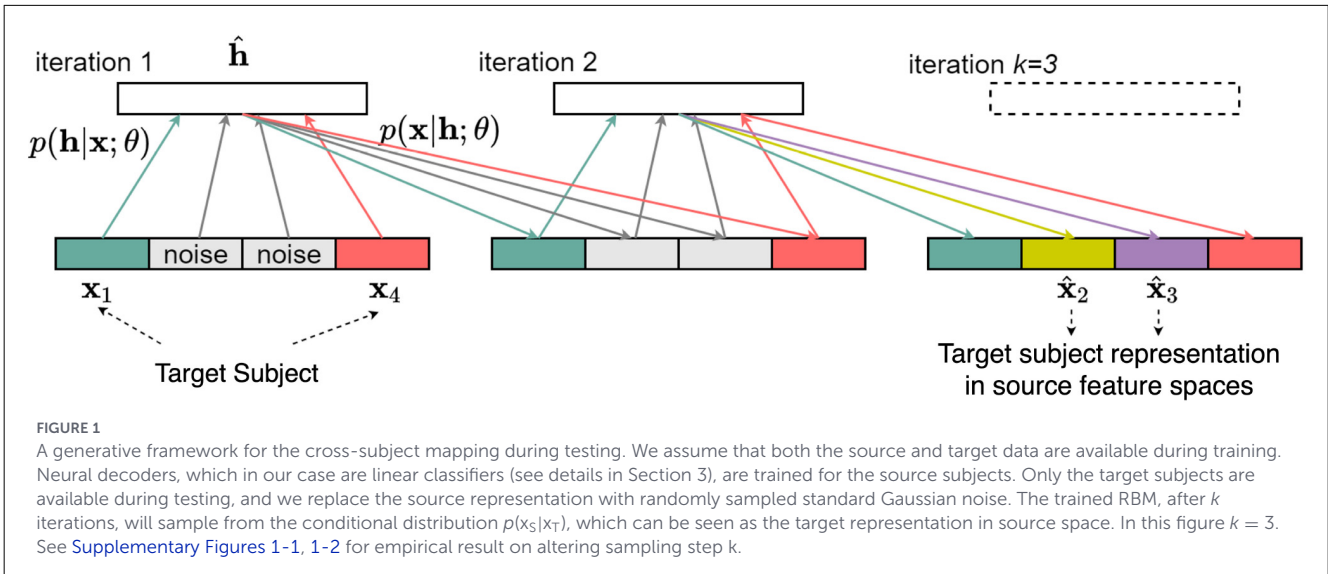
where the expectation is again taken w.r.t. the data-generating pdf $p_*(\mathbf{x})$. We note that $D_{\text{F}}(p_*, p) \geq 0$ with equality if and only if $p = p_*$ almost surely. Under mild regularity conditions, the Fisher divergence Equation 2 can be written as (Hyvärinen, 2005)

$$D_{\text{F}}(p_*, p) = \mathbb{E}_* [s_{\text{F}}(\mathbf{x}, \theta)] + c_* \quad (3)$$

where c_* is a term that does not depend on θ and $s_{\text{F}}(\mathbf{x}, \theta)$ is the Hyvärinen Score, defined as

$$s_{\text{F}}(\mathbf{x}, \theta) = \frac{1}{2} \|\nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta)\|^2 + \Delta_{\mathbf{x}} \log p(\mathbf{x}; \theta). \quad (4)$$

The result (Equation 3) enables to minimize the Fisher divergence over the space of parameters Θ by minimizing the empirical analog of the Hyvärinen Score $\bar{s}_{\text{F}}(\mathbf{x}_n, \theta)$. Let θ_* denote the parameter value that minimizes the Fisher divergence between $p(\mathbf{x}; \theta_*)$ and $p_*(\mathbf{x})$ between all model class candidates. By standard asymptotic analysis, it can be shown that the estimate $\hat{\theta}_{\text{F}} = \arg \min_{\theta} \bar{s}_{\text{F}}(\mathbf{x}_n, \theta)$ satisfies $\hat{\theta}_{\text{F}} \rightarrow \theta_*$ in probability as the



number of data points grows (we refer details to Hyvärinen, 2005 and references therein). This estimation procedure is known as score matching. It has been proved that score matching using the Langevin Monte Carlo method is equivalent to contrastive divergence in the limit of infinitesimal step size (Hyvärinen, 2007). Although this result implies that this variant of convergence divergence can retain the consistency on score matching, we note that this equivalence holds only for a particular Markov Chain Monte Carlo (MCMC) method. The actual performance of these two methods are different.

2.2.4 Gauss–Bernoulli restricted Boltzmann machines

An RBM is a bipartite undirected graphical model where only the links between visible units and hidden units are allowed. We focus on Gauss–Bernoulli RBM, which consists of continuous inputs $\mathbf{x} \in \mathbb{R}^D$ and binary hidden units $\mathbf{h} \in \{0, 1\}^H$ with the pdf

$$p(\mathbf{x}, \mathbf{h}; \theta) = \frac{e^{-E(\mathbf{x}, \mathbf{h}; \theta)}}{Z(\theta)}, \quad Z(\theta) = \sum_{\mathbf{h}} \int_{\mathbf{x}} e^{-E(\mathbf{x}, \mathbf{h}; \theta)} d\mathbf{x}, \quad (5)$$

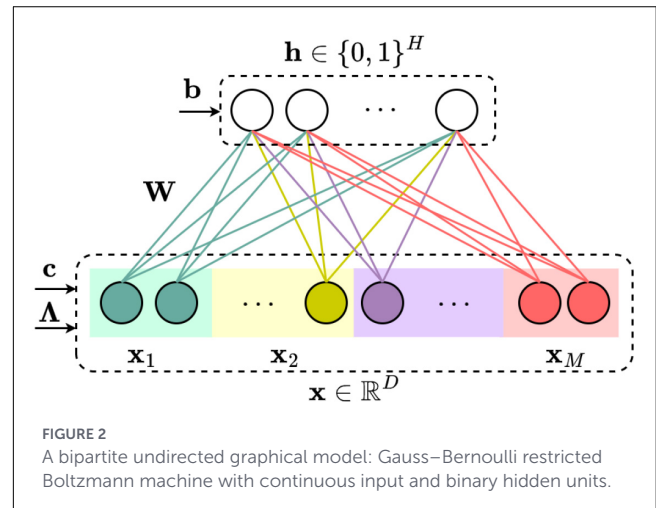
where the energy function $E(\mathbf{x}, \mathbf{h}; \theta)$ is given by

$$E(\mathbf{x}, \mathbf{h}; \theta) = \frac{1}{2}(\mathbf{x} - \mathbf{c})^\top \Lambda (\mathbf{x} - \mathbf{c}) - \mathbf{h}^\top \mathbf{W} \Lambda \mathbf{x} - \mathbf{b}^\top \mathbf{h}. \quad (6)$$

Here, $\mathbf{W} \in \mathbb{R}^{H \times D}$ is the matrix of weights connecting the hidden and input layer, $\mathbf{b} \in \mathbb{R}^H$ and $\mathbf{c} \in \mathbb{R}^D$ are the vectors of hidden and input layer biases, and $\Lambda = \text{dg}(\lambda)$ denotes the diagonal precision matrix of the inputs. All these parameters are freely learnable and they are denoted by θ in Equations 5, 6. We illustrate a Gauss–Bernoulli RBM model in Figure 2. It is easy to see that the conditional densities are given by

$$p(\mathbf{h}|\mathbf{x}; \theta) = \sigma(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b}), \quad (7)$$

$$p(\mathbf{x}|\mathbf{h}; \theta) = \mathcal{N}(\mathbf{W}^T \mathbf{h} + \mathbf{c}, \Lambda^{-1}). \quad (8)$$



The marginal density $p(\mathbf{x}; \theta)$ of the visible inputs can be also written in the energy-based form

$$p(\mathbf{x}; \theta) = \frac{e^{-\mathcal{F}(\mathbf{x}; \theta)}}{Z(\theta)}, \quad (9)$$

where $Z(\theta)$ is called the normalizing constant, and $\mathcal{F}(\mathbf{x}; \theta)$ is the free energy:

$$Z(\theta) = \int_{\mathbf{x}} e^{-\mathcal{F}(\mathbf{x}; \theta)} d\mathbf{x}, \quad \mathcal{F}(\mathbf{x}; \theta) = \frac{1}{2}(\mathbf{x} - \mathbf{c})^\top \Lambda (\mathbf{x} - \mathbf{c}) - \mathbf{1}_H^\top \boldsymbol{\gamma}, \quad (10)$$

with $\boldsymbol{\gamma} = \log(\mathbf{1}_H + \exp(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b}))$ denoting the element-wise Softplus function. Unlike Equation 6, the energy function (Equation 10) associated with the marginal $p(\mathbf{x}; \theta)$ is no longer linear in the free parameters θ .

A frequently encountered Gauss–Bernoulli RBM in the literature is the one associated with the conditional density $p(\mathbf{x}|\mathbf{h}; \theta) = \mathcal{N}(\mathbf{W}^T \mathbf{h} + \mathbf{c}, \mathbf{I}_D)$ and assumes unit variances for the input units. This is a special case of our model Equation 6 in

which we treat the variances of the inputs as learnable parameters, and all results and discussions in this paper can be applied in a straightforward manner to the special case by replacing Λ with the identity matrix.

2.2.5 Learning RBM via contrastive divergence

The negative log-likelihood of the parameters of the RBM, i.e., the logarithmic loss can be written as

$$\ell(\mathbf{x}; \boldsymbol{\theta}) \equiv -\log p(\mathbf{x}; \boldsymbol{\theta}) = \mathcal{F}(\mathbf{x}; \boldsymbol{\theta}) + \log Z(\boldsymbol{\theta}). \quad (11)$$

The gradient obtains a particularly interesting form:

$$\begin{aligned} -\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \mathcal{F}(\mathbf{x}; \boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) \\ &= \nabla_{\boldsymbol{\theta}} \mathcal{F}(\mathbf{x}; \boldsymbol{\theta}) - \mathbb{E} [\nabla_{\boldsymbol{\theta}} \mathcal{F}(\mathbf{x}; \boldsymbol{\theta})] \end{aligned} \quad (12)$$

where the expectation in the second term is taken w.r.t. the marginal density of the visible units given in Equation 9. Therefore, it is difficult to determine the gradient analytically. In order to make the computation tractable, this expectation is estimated using samples from $p(\mathbf{x}; \boldsymbol{\theta})$ which can be obtained by running a Markov chain with Gibbs sampling as the intermediate sampling operator. To speed up the sampling process, Hinton (2002) showed that it is not necessary to wait for the Markov chain to converge; instead, if the chain is initialized using training examples, reasonable learning performance might be obtained only after k Gibbs steps. In practice, $k = 1$ is commonly used during training (Hinton, 2002). However, this corresponds to the approximate minimization of the contrastive divergence (CD), which produces biased estimates of the model parameters (Carreira-Perpinan and Hinton, 2005). While obtaining an expected value may require a large k for a contrastive divergence to converge (van der Plas et al., 2023), in our application, we need to sample from the conditional distribution; therefore, we chose $k = 1$ during test time where we sample from the distribution of source conditioned on target. In the extended data, we empirically show that the choice of k only marginally changes the result and does not affect the conclusion of this work (see Supplementary Figures 1-1, 1-2 of extended data).

We see that an important implication of approximating MLE-based learning through contrastive divergence minimization is the lack of consistency guarantees. Specifically, minimizing the contrastive divergence is not guaranteed to converge to the data-generating parameter $\boldsymbol{\theta}_*$ that minimizes the KL divergence from $p(\mathbf{x}; \boldsymbol{\theta})$ to the data-generating pdf $p_*(\mathbf{x})$. The impediment can be traced back to the computation of the gradient of the logarithmic loss and the analytical intractability of the second term in Equation 12 which appears due to the intractability of the partition function as a normalizing constant in Equation 5.

2.2.6 Learning RBM via Fisher divergence

To overcome the issues associated with the lack of consistency guarantees, instead of aiming to minimize the

KL divergence through contrastive divergence approximation, we propose to minimize the Fisher divergence from the marginal density of the visible units $p(\mathbf{x}; \boldsymbol{\theta})$ to the data-generating distribution $p_*(\mathbf{x})$. To evaluate the Hyvärinen Score Equation 4 based on Equation 9, we derived the following result.

Proposition 1. *The Hyvärinen Score for the Gauss–Bernoulli RBM Equation 5 with energy function Equation 6 is given by*

$$s_F(\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} \|\Lambda(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x})\|^2 + \text{tr}(-\Lambda + \Lambda \mathbf{W}^T \text{dg}(\boldsymbol{\sigma}') \mathbf{W} \Lambda), \quad (13)$$

with $\boldsymbol{\sigma} = \sigma(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})$ and $\boldsymbol{\sigma}' = \sigma'(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})$, where σ and σ' respectively denote the element-wise Sigmoid operator and the corresponding first derivative.

Proof: Taking the derivative of the log-density $\log p(\mathbf{x}; \boldsymbol{\theta})$ w.r.t. \mathbf{x} , we obtain

$$\begin{aligned} \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta}) &= -\nabla_{\mathbf{x}} \mathcal{F}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\mathbf{x}} \log Z(\boldsymbol{\theta}) \\ &= \Lambda(\mathbf{c} - \mathbf{x}) + \sum_{h=1}^H \nabla_{\mathbf{x}} \mathcal{Y}(\mathbf{W}_h: \Lambda \mathbf{x} + b_h) \\ &= \Lambda(\mathbf{c} - \mathbf{x}) + \sum_{h=1}^H \sigma(\mathbf{W}_h: \Lambda \mathbf{x} + b_h) \Lambda \mathbf{W}_h^T \\ &= \Lambda(\mathbf{c} - \mathbf{x}) + \Lambda \mathbf{W}^T \sigma(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b}), \end{aligned}$$

which gives the first term in Equation 13. To obtain the second term, we first compute the Hessian matrix of the log-density $\log p(\mathbf{x}; \boldsymbol{\theta})$; we obtain:

$$\begin{aligned} \nabla_{\mathbf{x}\mathbf{x}}^2 \log p(\mathbf{x}; \boldsymbol{\theta}) &= \nabla_{\mathbf{x}} (\nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta})) \\ &= \nabla_{\mathbf{x}} (\Lambda(\mathbf{c} - \mathbf{x}) + \Lambda \mathbf{W}^T \sigma(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})) \\ &= -\Lambda + \nabla_{\mathbf{x}} \sigma(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b}) \mathbf{W} \Lambda \\ &= -\Lambda + \Lambda \mathbf{W}^T \nabla_{\mathbf{u}} \sigma(\mathbf{u}) \mathbf{W} \Lambda \\ &= -\Lambda + \Lambda \mathbf{W}^T \text{dg}(\boldsymbol{\sigma}'(\mathbf{u})) \mathbf{W} \Lambda, \end{aligned}$$

where $\mathbf{u} = \mathbf{W} \mathbf{x} + \mathbf{b}$. Plugging the Hessian into the Laplacian $\Delta_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \text{tr}(\nabla_{\mathbf{x}\mathbf{x}}^2 \log p(\mathbf{x}; \boldsymbol{\theta}))$ gives the second term in Equation 13, which completes the proof.

We observe that unlike the logarithmic loss in Equation 11, the Hyvärinen Score can be evaluated explicitly in terms of the parameters \mathbf{c} , \mathbf{b} and \mathbf{W} of the Gauss–Bernoulli RBM. Moreover, the calculation does not involve the partition function $Z(\boldsymbol{\theta})$. This simplifies the computation of the gradient w.r.t. the parameters of the RBM, which now can be computed by straightforward

application of matrix calculus yielding the following closed-form expressions:

$$\begin{aligned} \nabla_{\mathbf{c}} s_{\mathbf{F}}(\mathbf{x}, \boldsymbol{\theta}) &= \Lambda^2(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}), \\ \nabla_{\mathbf{b}} s_{\mathbf{F}}(\mathbf{x}, \boldsymbol{\theta}) &= \text{dg}(\boldsymbol{\sigma}') \mathbf{W} \Lambda^2(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}) + \text{dg}(\mathbf{W} \Lambda^2 \mathbf{W}^T) \boldsymbol{\sigma}'', \\ \nabla_{\mathbf{W}} s_{\mathbf{F}}(\mathbf{x}, \boldsymbol{\theta}) &= \text{dg}(\boldsymbol{\sigma}') \mathbf{W} \Lambda^2(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}) \mathbf{x}^T \Lambda \\ &\quad + \boldsymbol{\sigma}(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x})^T \Lambda^2 + \text{dg}(\mathbf{W} \Lambda^2 \mathbf{W}^T) \boldsymbol{\sigma}'' \mathbf{x}^T \Lambda \\ &\quad + 2 \text{dg}(\boldsymbol{\sigma}') \mathbf{W} \Lambda^2, \\ \nabla_{\lambda} s_{\mathbf{F}}(\mathbf{x}, \boldsymbol{\theta}) &= \text{dg}(\mathbf{x}) \mathbf{W}^T \text{dg}(\boldsymbol{\sigma}') \mathbf{W} \Lambda^2(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}) \\ &\quad + \text{dg}(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}) \Lambda(\mathbf{W}^T \boldsymbol{\sigma} + \mathbf{c} - \mathbf{x}) \\ &\quad + 2 \text{dg}(\mathbf{W}^T \text{dg}(\boldsymbol{\sigma}') \mathbf{W}) \lambda + \mathbf{1}_N \\ &\quad + \text{dg}(\mathbf{x}) \mathbf{W}^T \text{dg}(\mathbf{W} \Lambda^2 \mathbf{W}^T) \boldsymbol{\sigma}'' \end{aligned}$$

with $\boldsymbol{\sigma} = \boldsymbol{\sigma}(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})$, $\boldsymbol{\sigma}' = \boldsymbol{\sigma}'(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})$ (as in Proposition 1), whereas $\boldsymbol{\sigma}'' = \boldsymbol{\sigma}''(\mathbf{W} \Lambda \mathbf{x} + \mathbf{b})$; also, recall that $\Lambda = \text{dg}(\lambda)$.¹

It is evident that, as opposed to the gradient of the logarithmic loss (Equation 12), the gradient of the Hyvärinen Score can be computed explicitly w.r.t. the parameters of the Gauss–Bernoulli RBM, producing closed-form expressions that can be used directly for training the parameters of the RBM. Indeed, let $\hat{\boldsymbol{\theta}}_{\mathbf{F}}^{(n)}$ denote the parameter estimate at each step n . The new parameter estimate can be obtained through the following update rule (repeated until convergence):

$$\hat{\boldsymbol{\theta}}_{\mathbf{F}}^{(n+1)} = \hat{\boldsymbol{\theta}}_{\mathbf{F}}^{(n)} - \eta \frac{1}{B} \sum_b \nabla_{\boldsymbol{\theta}} s_{\mathbf{F}}(\mathbf{x}_b, \hat{\boldsymbol{\theta}}_{\mathbf{F}}^{(n)}),$$

where η is the learning rate, B is the size of the minibatch of randomly chosen data points \mathbf{x}_b , and $b = 1, \dots, B$ at step n .

2.3 Cross-subject mapping algorithm

Recall from Section 2.1 that in cross-subject mapping, the goal is to obtain source feature representation(s) \mathbf{x}_S from the target feature vector(s) \mathbf{x}_T . We will elaborate on how to use an RBM and the Gibbs sampler to sample such representations. We first parameterize the generative model $p(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$ of all feature vectors $\mathbf{x} = (\mathbf{x}_S, \mathbf{x}_T)$ and hidden variables \mathbf{h} using the Gauss–Bernoulli RBM described in Section 2.2.4, see also Figure 2. After learning the parameters of the Gauss–Bernoulli model we infer $\mathbf{x}_S = (\mathbf{x}_i)_{i \in \mathcal{M}_S}$ from $\mathbf{x}_T = (\mathbf{x}_j)_{j \in \mathcal{M}_T}$ as follows. First, we initialize the vector $\hat{\mathbf{x}}$ by the features of the target subjects \mathbf{x}_m , $m \in \mathcal{M}_T$ and random noise (e.g., with standard normal variables). Then:

1. generate $\hat{\mathbf{h}} \sim p(\mathbf{h}|\hat{\mathbf{x}}; \boldsymbol{\theta})$ via Equation 7;
2. using $\hat{\mathbf{h}}$ generate $\hat{\mathbf{x}} \sim p(\mathbf{x}|\hat{\mathbf{h}}; \boldsymbol{\theta})$ via Equation 8.

We obtain the final estimate after repeating the above two steps $k \geq 1$ times; Figure 1 illustrates an example with $k = 3$. This gives the source feature space representations of the target feature vectors and they can be further processed using algorithms trained on source data. Alternatively, in the final step, we can skip sampling

from $p(\mathbf{x}|\hat{\mathbf{h}}; \boldsymbol{\theta})$ and we can also infer \mathbf{x}_S as $\hat{\mathbf{x}}_S = (\hat{\mathbf{x}}_j)_{j \in \mathcal{M}_S} = \max_{\mathbf{x}_S} p(\mathbf{x}|\hat{\mathbf{h}}; \boldsymbol{\theta})$. For simplicity and without loss of generality, we have assumed that target and source subjects \mathcal{M}_T and \mathcal{M}_S satisfy $\mathcal{M} = \mathcal{M}_T \cup \mathcal{M}_S$ and $\mathcal{M}_T \cap \mathcal{M}_S = \emptyset$. This allows us to skip a tedious step in the algorithm and avoid the marginalization of the joint density $p(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$ over subjects that are neither targets nor sources.

3 Evaluation and results

Next, we present the results from the evaluations. First, in Section 3.1 we describe the experimental protocol and the acquired data. In Section 3.2 we discuss the evaluation methodology, including evaluation scenarios. In Section 3.3 we present the main findings and observations.

3.1 Experiment, data and features

3.1.1 Protocol

We study a comprehensive flight motor program for hawk moths. We will describe the experimental protocol and related procedures only briefly here; the interested reader is referred to Putney et al. (2021) where the data set was first published for more details. The subjects, i.e., the moths are tethered inside a three-sided box formed by computer monitors displaying the visual stimuli. Each stimulus is represented by sinusoidal gratings with a spatial frequency of 20° per cycle on 3D spheres projected on the monitors. The spheres drift at a constant velocity of 100° per second, corresponding to a temporal frequency of five cycles/s. Moreover, the spheres also drift in opposite directions about the three axes of rotation which result in six different visual stimuli also known as pitch (up, down), roll (left, right), and yaw (left, right) (Putney et al., 2021).

The moth responds to each of the six discrete stimuli by producing turning effort as assessed with a six degree-of-freedom force/torque transducer. The 10 primary muscles that control the flying motion of the moth are wired and enable spike-resolved EMG signals to be recorded during tethered flight. These key muscles include the flight power muscles (dorsolongitudinal (DLM) and dorsoventral (DVM) muscles), as well as the steering muscles [third axillary (3AX), basalar (BA), and subalar (SA) muscles] on both the left and the right side of the thorax. The EMG recordings are used to extract the timings of the motor unit spikes in each of the muscles that serve as control commands by means of which the nervous system guides the motion of the moth in response to the different visual stimuli; more details can be found in Putney et al. (2019) and Putney et al. (2021). Taken together, this dataset is unusual in its near-complete recording of all the spikes the animal can use to control its wings and so is an ideal point of convergence to test for the decodability of stimulus conditions.

The objective is to decode the visual stimulus from the comprehensive motor program recordings, i.e., the spike trains. The subject-specific formulation of the problem where the neural decoder (classifier) is both trained and tested on the same

¹ We omit the detailed derivation of the gradients for brevity, and we note that the gradients can be derived through straightforward application of matrix calculus.

subject was analyzed in Putney et al. (2021). Here, we study the performance of the neural decoder in cross-subject settings, where the test data originates from the target subjects whereas the classifier is trained on source subject data.

The dataset is collected from nine subjects of either sex with over 20 seconds of recording sessions for each visual stimulus. Each session is segmented into *wing strokes*, i.e., trials (Putney et al., 2021). The typical duration of a wing stroke is between 50 and 70 ms yielding an average number of trials of $\approx 2,500$ per moth. It should be noted that some moths in the dataset are missing the recordings from some of their muscles (either one or at most two) due to failures in the recording procedure. Nevertheless, as demonstrated in Putney et al. (2021), the absence of some (one or two) muscles does not have a significant effect on decoding performance; in fact, as shown in Putney et al. (2021), high decoding accuracy (higher than 90%) can be achieved even with half of the available muscles due to the completeness of the motor program.

3.1.2 Feature extraction

Before we delve into more details with respect to the cross-subject neural decoder, we briefly describe our methodology for constructing feature representations from spike trains proposed in Putney et al. (2021). Since the spike trains are given by variable-length vectors of spike timings, we consider Gaussian kernels, a strategy commonly used in neuroscience, to interpolate the spike trains. The Gaussian kernel is given by:

$$x(t) = \sum_n \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right), \quad 0 \leq t \leq \tau,$$

where t_n denotes the timing of the n -th spike collected from an arbitrary muscle, trial, and moth, τ denotes the wing stroke cut-off threshold (we only consider spikes that satisfy $t_n \leq \tau$), and σ is the Gaussian kernel bandwidth. The goal is to obtain a smooth multivariate time-series representation of the spike trains in which the spike timing information is conveyed by centering one kernel at each spike and summing the kernels; these yields feature vectors of fixed dimension $\tau \cdot \nu_S$ where ν_S is the sampling frequency. For consistency, the muscles whose recordings are missing are filled with zero vectors of the same length as above. We then flatten the interpolated time series across muscles to obtain one large feature vector. Finally, we apply PCA and retain only the first P largest modes; this is our final representation \mathbf{x}_m of the neural activity of subject m with dimension $D_m = P$. Figure 3 depicts the feature spaces of each moth and all six visual conditions (tasks) across the first two strongest principal components, i.e., modes after performing PCA. The diagrams also show the confidence ellipses corresponding to one standard deviation for each of the conditions. It can be clearly observed that the data demonstrates strong separability properties even in the first two dimensions of the PCA-based feature space; adding more features (i.e., PCA modes) only increases this separability in the higher-dimensional feature space for each moth as we have observed in our past work (Putney et al., 2021); for more details, we advise the interested reader to refer to Putney et al. (2021) where the feature extraction procedure was first proposed and its performance thoroughly analyzed.

3.2 Scenarios and methodology

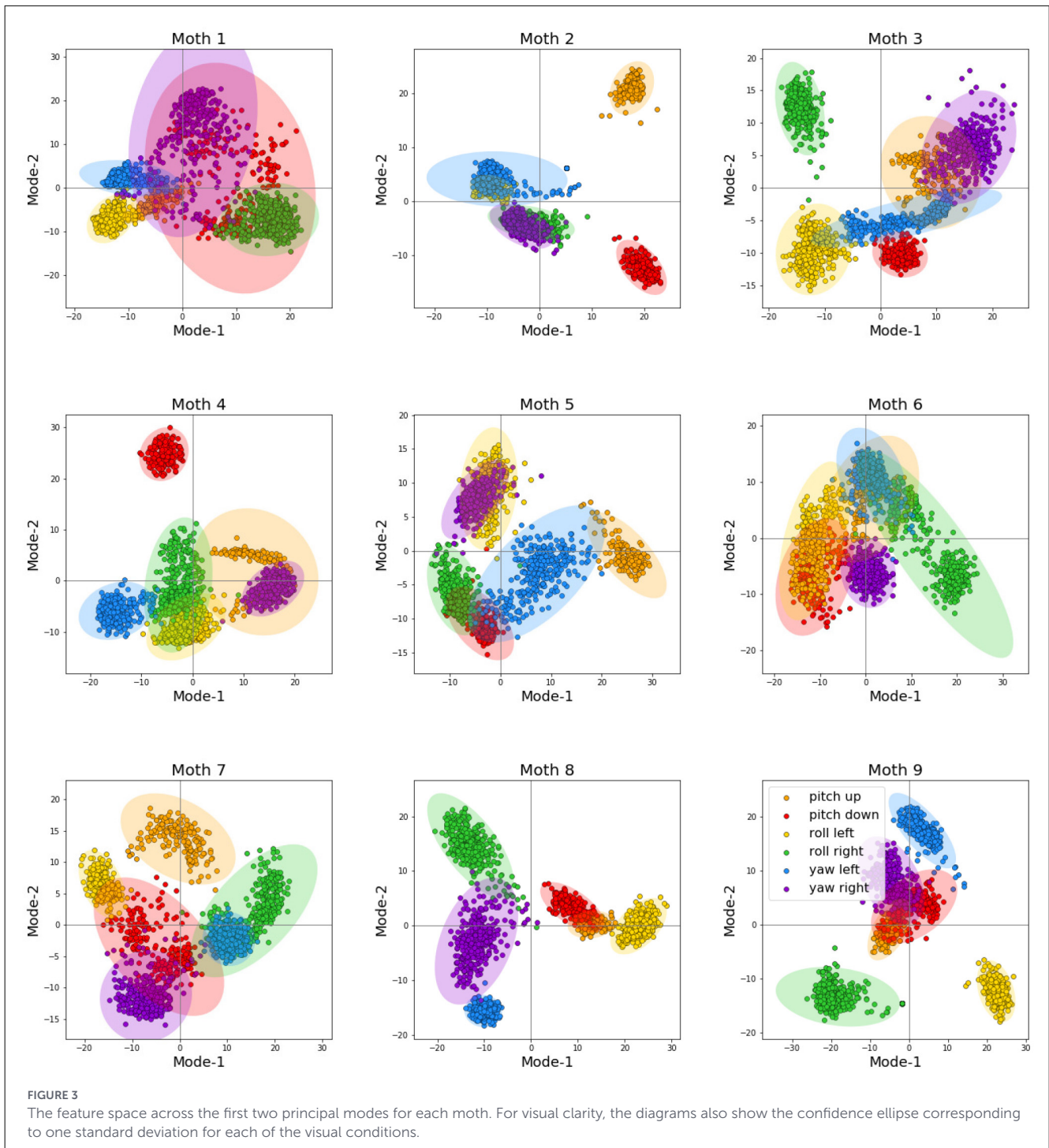
The moth population set is $\mathcal{M} = \{1, \dots, 9\}$. We evaluate the performance of the cross-subject neural decoder in the following scenarios:

- *Scenario I* (Figure 4, left). The source and target sets are $\mathcal{M}_S = \{m\}$ and $\mathcal{M}_T = \mathcal{M} \setminus \{m\}$, respectively. That is, we select a single subject $m \in \mathcal{M}$ as the source and map the features of *all* remaining subjects $i \neq m$ onto the feature space of the source subject m . The obtained representation is then decoded using a linear classifier trained on subject m data.
- *Scenario II* (Figure 4, right). The target and source index sets are $\mathcal{M}_T = \{m\}$ and $\mathcal{M}_S = \mathcal{M} \setminus \{m\}$, respectively. In other words, we select a single subject $m \in \mathcal{M}$ as the target and map the corresponding feature vector onto the feature spaces of *all* remaining subjects $j \neq m$. Similarly, as in scenario I above, the obtained representations of the target feature vector are subsequently decoded using the subject-specific linear classifiers trained on the source subjects $j \neq m$ individually.

In both of these scenarios, we evaluate the performance of the source classifier (trained purely on source data) on the transferred target features through an RBM model with both standard contrastive divergence minimization and Fisher divergence minimization; we use RBM-CD and RBM-FD to denote these two RBM models, with CD standing for contrastive divergence (see Section 2.2.5) and FD standing for Fisher divergence (see Section 2.2.6). We compare the performance with three benchmarks:

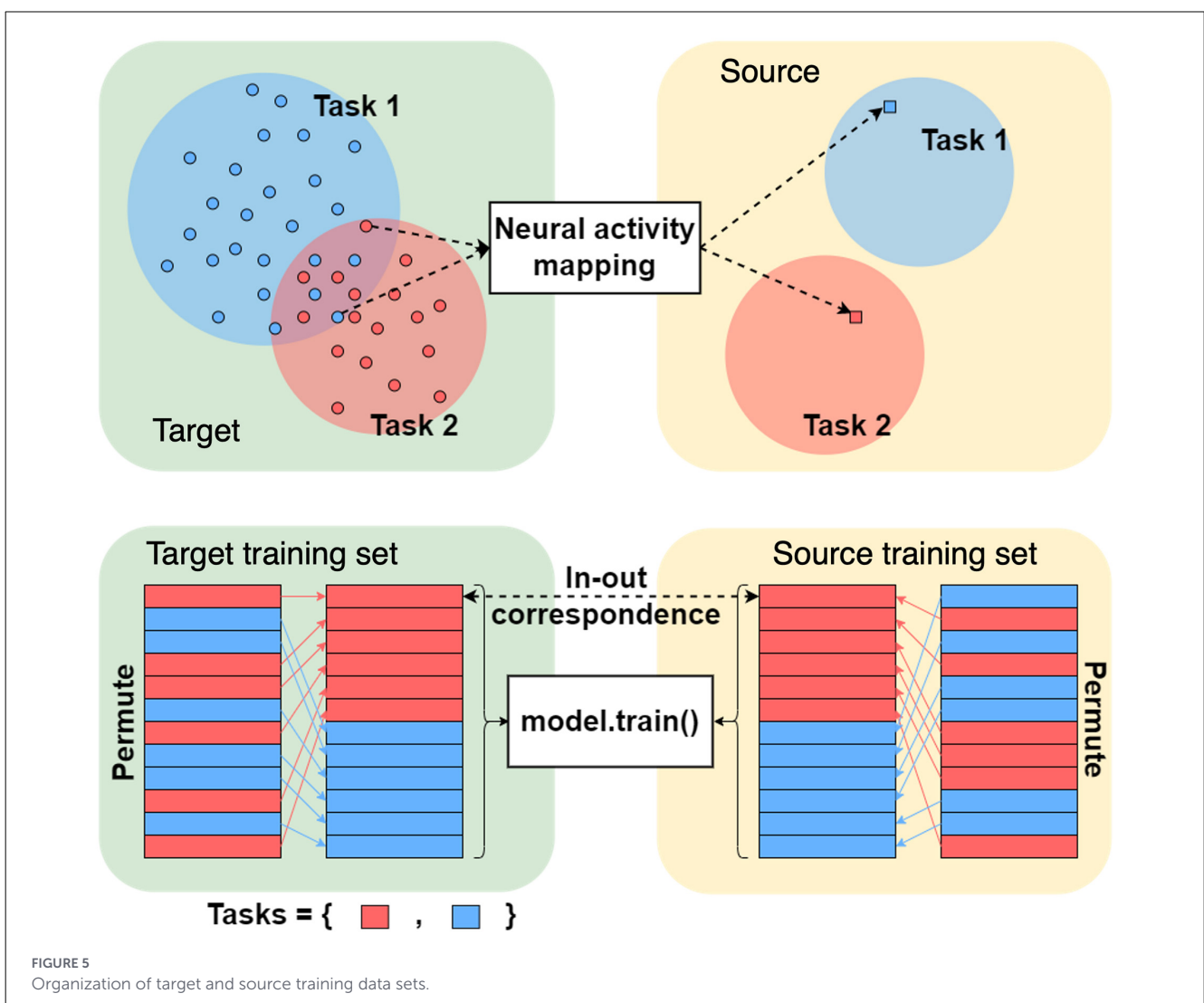
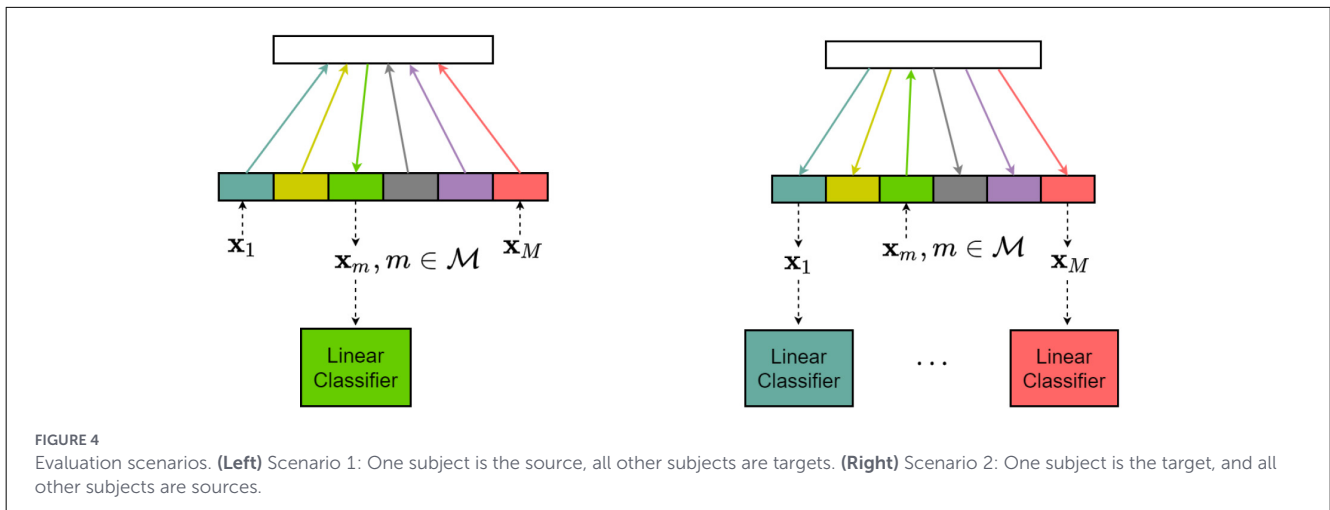
1. Subject-specific neural decoding, when the classifier is trained *and* tested on purely source data. Equivalently, the data used to train the classifier and the test data respectively come from the same individual subject. In this case, the performance of the neural decoder may be thought of as an upper bound because it represents the best-case scenario where the classifier has access to the most relevant information about feature patterns of the test data.
2. Cross-subject neural decoding with *no* transfer, when the target data is directly decoded using the source neural decoder without transferring the target data into source feature space. As we discussed in Section 1, in this case, the performance of the neural decoder is usually close to a random guess, since the neural decoder does not take into account the differences between neural activity patterns between the test (target) and the train (source) data.
3. Cross-subject neural decoding using HiWA, as described in Section 1, is a supervised state-of-the-art linear alignment technique that exploits the cluster structure to solve a hierarchical optimal transport optimization (Lee et al., 2019). This method is widely applied in neural settings (Dabagia et al., 2023).

The training and testing data for the target and source data sets in both scenarios are formed by splitting the original data sets of each moth into training and testing data subsets randomly according to $\omega \in (0, 1)$, which denotes



the ratio between the number of training samples and the total number of trials. To characterize the performance of the cross-subject neural decoder statistically, the random train/test data splitting procedure is repeated 100 times; the entire model including the RBM is then re-trained using the new training data, and the test performance is recorded. We also perform a five-fold cross-validation analysis to compare Fisher divergence and contrastive divergence RBM. The result of this cross-validation experiment is provided in [Supplementary Figures 6-1, 7-1](#).

In cross-subject mapping, we are aiming to find a source space feature representation of the target feature vector; the mapping should be consistent with the task that the target features are encoding. In other words, the mapping should be such that the transferred representation of the target feature vector should be in the region of the source feature space that corresponds to the task that the target is performing, as illustrated in the top row in [Figure 5](#). In essence, when sampling from the joint pdf $p(\mathbf{x}, \mathbf{h})$, we wish to sample from a stimulus-dependent distribution and this should be guided by the task



information the target feature vector is encoding. Hence, the trials in the training data set should be organized such that there is an input-output correspondence with respect to the stimulus. As we are unable to establish correspondence between

individual training trials, we assign the correspondences at random. The procedure is schematically depicted in the bottom row in Figure 5. Namely, for each target feature vector from the training set, we randomly choose a source feature vector from

TABLE 1 Experimental parameters.

Parameter	Value
Data split ratio (ω)	0.5
Wing stroke cut-off (τ)	60 ms
Kernel bandwidth (σ)	2.5 ms
Number of principal modes ($D_m = P$)	10
Number of units in hidden layer (H)	15
Optimizer for training RBM model	Adam (Kingma and Ba, 2014)
Learning rate	0.005
Minibatch size	150
Training epochs	350

the same stimulus class and declare the pair to be an input-output pair.

3.3 Results

The hyperparameter used in this work is provided in Table 1.

In both scenarios, we use the Linear Discriminant Analysis (LDA) for classification, which is also applied in Putney et al. (2021) and has been shown to perform exceptionally high decoding accuracy. LDA is trained only with the source data because under our general utility described in Section 2.1, motor intentions are collected only for the source subject. We also note that only one RBM is trained for both scenarios, reflecting the efficiency of our architectures under multi-subjects cross-subject mapping regime.

The results are shown in Figures 6, 7. We observe that in both scenarios the performance of the cross-subject neural decoder is bounded between the performances of the first and second benchmarks, with the performance of the subject-specific neural decoder being the upper bound and the performance of the cross-subject neural decoding with no transfer being the lower bound. Note that in Figure 6 the lower bound varies around 0.16, which corresponds to random choice decoding in our case (as there are six stimuli in the experiment, see Section 3.1). The performance of the cross-subject neural decoder is similar in Scenario II but we omit to show it in Figure 7 to avoid clutter.

We conclude that the performance of the cross-subject neural decoder without transferring the target features to the source is very poor, and produces a poor lower bound, which is also intuitively expected. This can be most easily seen by inspecting Figure 3 which depicts the feature spaces the first two modes for each moth. Even though each moth individually exhibits a high degree of class separability (which ultimately results in very reliable subject-specific decoding performance as demonstrated in Figure 6), there is very little alignment between the geometric distribution of the classes/tasks (i.e., visual conditions) in the space spanned by the first two modes across different moths. In fact, the task-specific representations seem to occupy arbitrary segments of the feature

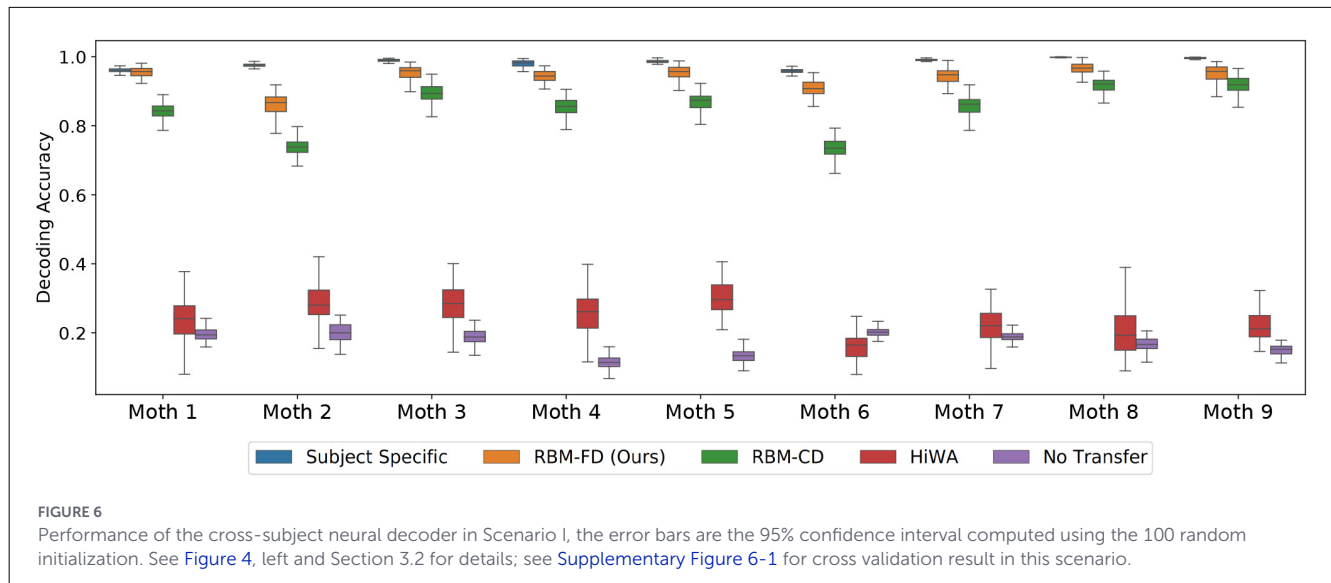
space across the first two dimensions and no discernible pattern can be directly observed; adding even more modes/features which are required to achieve high subject-specific separability, only exacerbates the differences of the feature spaces across moths. As a result, directly decoding any target moth using a neural decoder trained over a different, source moth results in poor performance as reported in Figure 6.

While HiWA does improve upon the lower bound (cross-subject decoding with no transfer), RBM methods consistently outperform HiWA with higher downstream classification accuracy and smaller variances under different random splitting and initialization in both scenarios. With the complex nature of neural data across subjects, there is no guarantee that different subjects can be mapped through a linear alignment, which could be the reason for HiWA's limited performance in this dataset. RBM learns a non-linear mapping function that adapts the target to the source feature space where the decoder was trained with few assumptions; in a highly complex setting with multiple subjects and high dimensionality, which is a common theme in neural datasets, our RBM methods could provide a more adaptable alternative to the widely applied alignment methods.

The goal of cross-subject transfer learning is to learn a mapping function that takes target features and adapts them to the spaces of the source where the decoder was trained. We present the RBM-adapted feature space in Figure 8 which shows the distribution of the target test points in Scenario I after applying the FD-RBM cross-subject transfer with trained RBM model to the corresponding source feature space. Note that the test points are the features coming from a diverse set of targets, namely all remaining (eight) moths; hence, the illustrative result shown in Figure 8, in addition to the decoding results presented in Figure 6, clearly demonstrate that the FD-RBM model has successfully learned a non-linear transformation that takes a task-specific feature representation from an arbitrary target and maps it into the adequate task-specific region of the source feature space.

By comparing the results in Figure 6 with the results in Figure 7, we also observe that the performance of the neural decoder in Scenario I outperforms the neural decoders in Scenario II. This is an intuitively expected result, as in Scenario II, the size of the joint source feature vector \mathbf{x}_S is $M - 1$ times larger than the target feature vector \mathbf{x}_T ; that is, in Scenario II we are jointly obtaining the representation of a single target in 8 different source feature spaces. The opposite reasoning applies to Scenario I. Hence, the drop in the performance from Figures 6 to 7 is expected. Furthermore, for a given population of subjects indexed in \mathcal{M} , we can view these two scenarios as the two extreme cases that put the upper (Scenario I) and lower (Scenario II) bounds on the performance.

We also observe that the performance of the cross-subject neural decoder with RBM-FD transfer outperforms the performance of the same cross-subject decoder with RBM-CD transfer. This is an interesting finding, further indicating that in the case of RBMs, the training based on Fisher divergence minimization yields better results in comparison with the more conventional approach based on Maximum Likelihood. One potential explanation for this behavior is that the gradient for Fisher divergence can be explicitly computed while contrastive



divergence requires an expectation through sampling to compute its gradient. While traditional contrastive divergence can be improved through persistent contrastive divergence (Tieleman, 2008), the explicit gradient form could lead to more stable training of the RBM. This result is also consistent with our findings on popular public datasets such as the MNIST where we used RBM for applications such as compression and reconstruction and where we observed that an RBM trained via Fisher divergence minimization yields higher-quality image reconstructions. In addition to the reliability improvement, we note that training an RBM-FD is less computationally demanding as opposed to RBM-CD which requires Gibbs sampling even during training to obtain estimates for the gradients. However, it should be noted that while the RBM-FD model in general tends to outperform the RBM-CD model, the behavior is ultimately determined by the values of the free parameters, and in the case of the parameters we have selected (listed in the beginning of this section), the above observations are valid.

4 Discussion

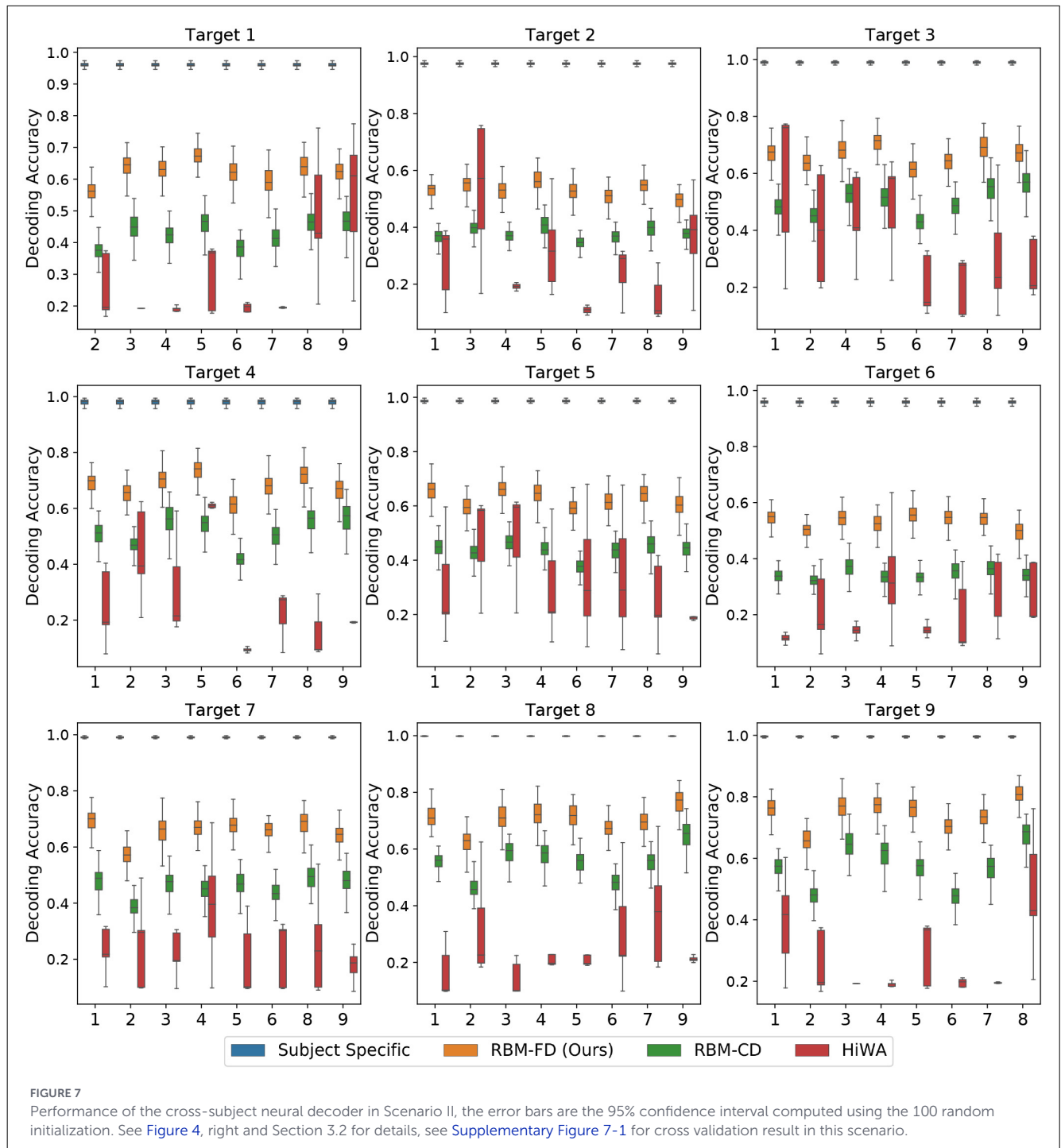
Designing reliable, robust, and computationally efficient solutions for cross-subject mapping and data integration is a major challenge in neuroscience. However, the demand for such algorithms, especially outside experimental settings where subject-specific neural decoders cannot be reliably trained, has grown significantly (CTRL-labs at Reality Labs et al., 2024). In this paper, we propose a novel framework for learning the joint distribution of target and source feature representations across subjects using an undirected graphical model and the flexible, non-linear generative RBM. To enhance efficiency, we introduce an alternative RBM training method that minimizes Fisher divergence, enabling closed-form gradient computation and eliminating the need for Gibbs sampling. We evaluated our approach against HiWA, a recent distribution alignment method, on a neural decoding task using

experimental data from nine hawk moths exposed to six visual stimuli. Our method demonstrated substantial improvement over HiWA in two transfer scenarios, confirming its viability. These flexible, assumption-light approaches show promise in generalizing complex neural data across individuals, tuning neural interfaces to subject-specific features, and leveraging multi-subject data when experiments are time-limited or incomplete.

4.1 Application of restricted Boltzmann machine in neuroscience

Restricted Boltzmann machines are widely applied as a tool to study neuroscience problems. The bipartite network structure has long attracted neural science interests to study the structural connection and functional activities in the brain (Hjelm et al., 2014). Recently, RBM has been applied to study the structural connectivity of zebrafish larvae and identified coarse-grained neural grouping (van der Plas et al., 2023). Combining a recurrent element with the RBM for temporal analysis, a variant of RBM, the recurrent temporal RBM was shown to capture the temporal dynamic of the whole brain activities of zebrafish larvae (Monnens et al., 2024).

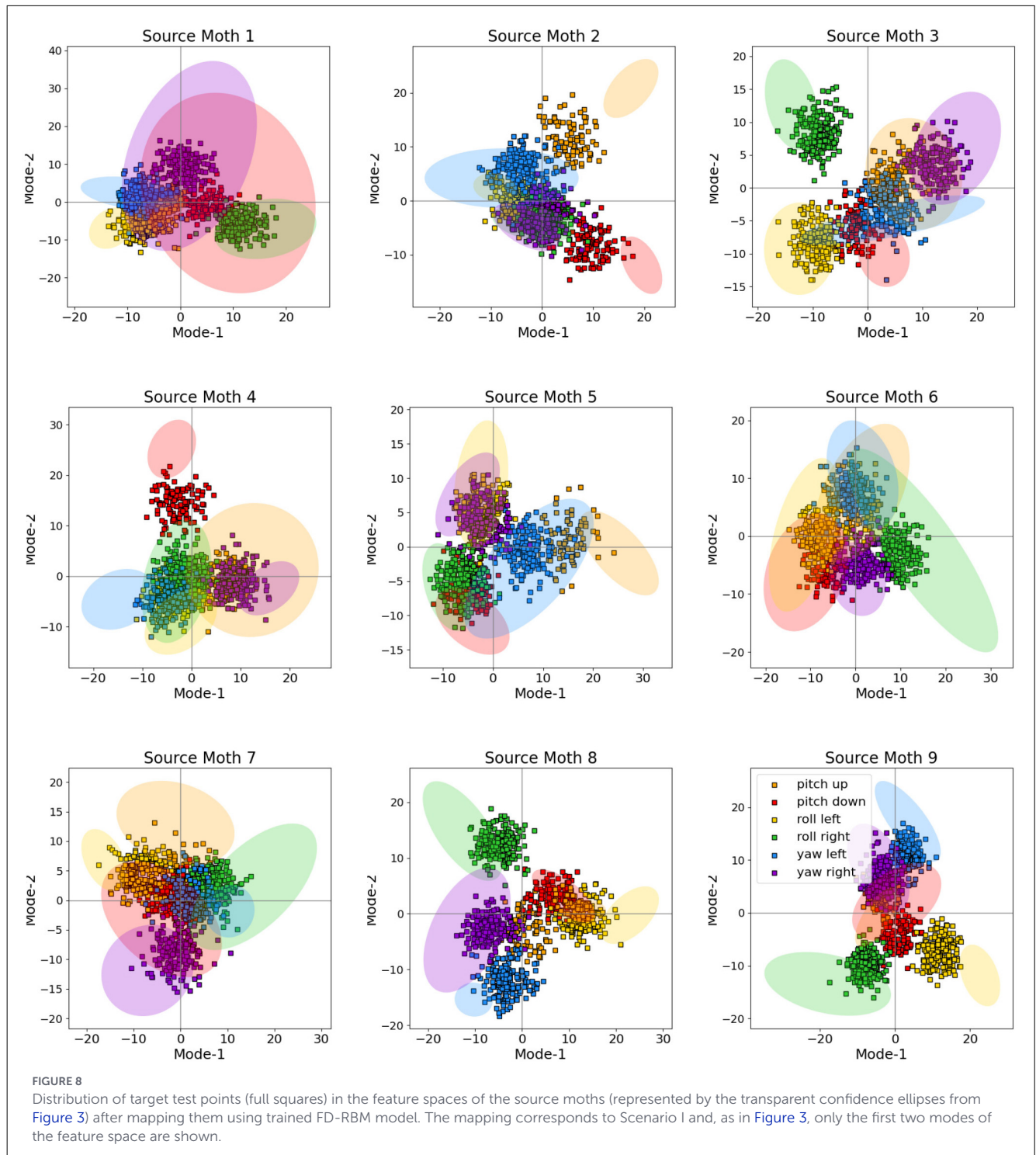
Our work considers a different application in neural science. We leverage the RBM's undirected graph structure to learn joint distributions of neural activities and then utilize this joint distribution to map multiple subjects for independent downstream tasks. Cross-subject mapping of EEG has been extensively explored in domains such as emotion detection (Li et al., 2018; Chen et al., 2021; Zhang et al., 2023). However, existing approaches are typically optimized for a fixed decoding objective and rely on task-specific classifiers, limiting their generalizability across paradigms. In contrast, the proposed generative framework aims to capture shared structure in neural activity itself, rather than features tailored to a particular behavioral or cognitive label.



This perspective is especially relevant for Brain—Computer Interface (BCI) applications, where each subject traditionally requires an extensive calibration phase prior to online operation, and where inter-subject variability in neural signals remains a major obstacle to decoder generalization (Fahimi et al., 2019; Zhong et al., 2024; Wang et al., 2024). While recent task-agnostic approaches have begun to address this issue (Wang et al., 2024), they often depend on large-scale feature extractors that increase computational complexity and reduce experimental flexibility. RBMs provide a comparatively simple

and interpretable generative model, with efficient sampling and a natural mechanism for learning shared structure across subjects, making them well suited for neuroscience applications where computational resources and experimental constraints are critical.

Finally, many neuroscience datasets exhibit strong temporal dependencies arising from ongoing neural dynamics and experimental structure. Extending the proposed framework to explicitly account for temporal organization, for example, through time-series RBM variants (Kuremoto et al., 2014; Tian and Xu,



2021), represents a natural and biologically relevant direction for future work.

In general, with its special bipartite network construction and undirected probabilistic graphical model structure, RBM is a highly flexible tool for further exploring neural science, understanding brain structures and dynamical neural activities, or applying it to cross-subject mapping.

4.2 Efficient and versatile cross-subject mapping

Many cross-subject mapping methods have been proposed to improve cross-subject mapping methods to be more accurate, robust, and versatile. One way to map between subjects is through alignment, which aligns the neural representation between two

subjects directly after simple preprocessing such as PCA (Lee et al., 2019; Gallego et al., 2020); one can also map the neural representation to a latent manifold before learning a mapping between the source and the target (Degenhart et al., 2020; Dabagia et al., 2023; Herrero-Vidal et al., 2021). More recently, generative algorithms emerging from deep learning were also applied to cross-subject mapping. For example, a Generative Adversarial Network was applied to transfer human activities between two subjects (Soleimani and Nazerfard, 2021) while a conditional VAE was applied to transfer latent neural representation between monkeys (Angelichinoski et al., 2020b).

However, limited studies focus on improving the general efficiency of cross-subject mapping methods, especially when multiple subjects are present in a complete dataset. Most of the cross-subject mapping methods deal with two subjects where at most two transformations are required to enable cross-subject mapping (Lee et al., 2019; Herrero-Vidal et al., 2021). As the number of subjects grows, the number of models needed to transfer between all subjects grows exponentially. To reduce the number of mappings when multiple subjects exist, researchers have developed a feature transfer that needs to be learned in a supervised manner (Zhao et al., 2021). Our work approaches a similar multi-subject/multi-target mapping problem from a probabilistic standpoint: by using RBM and estimating the joint distribution of all subjects at the same time, we learn the conditional probability that allows the mapping of an arbitrary number of targets to an arbitrary number of sources. This process improves the efficiency of cross-subject mapping with no required labels as we directly learn the distribution of features without a decoder.

Cross-subject mapping in a generative, probabilistic way through RBM also allows the RBM to be more versatile to apply in various settings. A common drawback of current methods of cross-subject mapping requires the availability of the source subject's data during test time. For example, the alignment-based algorithms can only correctly compute the alignment matrix when both the source and test are available (Lee et al., 2019; Degenhart et al., 2020). However, the RBM allows us to transfer different target subjects to the space of the source subject without additional knowledge regarding the source beyond its training data. This is enabled by the probabilistic modeling of RBM. With the joint distribution of the subjects represented by the RBM, conditional distribution can be easily computed for unknown subjects, allowing cross-subject mapping when a subject is not available. The probabilistic framework is also applied to align multiple subjects' neural responses for training a pooled decoder to classify odor stimuli (Herrero-Vidal et al., 2021).

4.3 Limitations

While the RBM architecture with Fisher divergence proposed in this study greatly improves the efficiency of cross-subject mapping over many subjects, limitations remain from several perspectives: its asymmetrical transformation, limited extensions to quantify similarities between subjects, limited generalizability with confounders between subjects, and its lack of structure in the hidden layer that removes interpretability from this architecture.

Many cross-subject mapping algorithms are symmetrical, meaning that transforming one subject to another subject can be directly applied to transform in the opposite directions with limited changes to the transformation function. Usually, the symmetric methods map two subjects onto a shared embedding space, which is widely applied in neural alignment methods such as (Harvey et al., 2024; Kriegeskorte et al., 2008). The RBM cross-subject mapping method is directional: the transformation from one subject to all of the other subjects learned by an RBM is characterized through complex neural networks that are non-linear. Directional transformations inevitably will make the transformation process less efficient, but the capability of mapping between multiple subjects at once compensates for the lack of symmetry in the transformation.

Another drawback of our RBM methods is the lack of a clear distance or similarity metric to quantify similarities between subjects. Many alignment-based methods, such as Harvey et al. (2024), extract similarity information once alignment is complete (see Klabunde et al., 2023 for a review). While RBM can estimate joint distributions and compute information-based metrics like conditional entropy, this metric is asymmetric and unsuitable for quantifying similarity. Mutual information, a symmetric measure, requires marginal distributions that our current architecture cannot compute. Modifying RBM to enable this is a potential future direction, but beyond the scope of this study.

The proposed framework is also limited in its reliance on training data drawn from the same underlying contextual distribution across subjects. While the RBM can, in principle, learn arbitrary joint distributions provided during training, it cannot reliably generalize to target data generated under contexts outside the training distribution. In such cases, increasing training data alone is insufficient to overcome distributional mismatch, as the limitation arises from a lack of shared structure rather than sample scarcity. Such limitation, in fact, is common for most cross-subject mapping methods. From a modeling perspective, extending the framework to settings with context mismatch would require explicitly modeling contextual variables, requiring modeling neural activity and context variable jointly through RBM; it is also possible to introduce mechanisms to detect and adapt to distributional shift such as distributionally robust methods through data augmentation (Hasan et al., 2025), remains an open direction for future work.

A related concern is that the use of random pairing between source and target samples assumes approximate exchangeability of neural activity within the same context. While this approximation is sufficient in the present setting, structured dependencies arising from latent sub-contextual factors may violate this assumption and reduce mapping fidelity. More generally, successful cross-subject mapping requires sufficient overlap in the support of neural activity distributions across subjects. When target-subject activity occupies regions of neural state space that are poorly represented or absent in the source data, the RBM is unable to reliably extrapolate, leading to systematic reconstruction errors and degraded decoding performance. Addressing such confounds may require explicit modeling or alignment of latent factors prior to joint distribution learning, representing an important extension of the current approach. This direction of research is commonly

studied in causal inference (Shalit et al., 2017), and the technique developed there can be directly applied prior to the application of RBM.

An additional limitation of the proposed framework is that decoder transfer implicitly assumes that the relationship between neural activity and behavioral variables is conserved across subjects up to a representational transformation. If subjects differ substantially in how neural activity maps to behavior—due to physiological differences (e.g., different motor program), or task strategy (e.g., different species)—then accurate neural alignment alone may not be sufficient for successful decoder transfer. From a neuroscience perspective, this limitation motivates experimental designs to stimulate behavioral variability with representation transformation, allowing such transformation to be recorded and modeled in future methods.

The lack of interpretability is another limitation of our current architecture. Cross-subject mapping methods that utilize algebraic solutions can identify how each feature of the target subject is mapped to the space of the source subject. For example, the Procrustes Alignment method computes a simple transformation matrix that maps the target data onto the source space (Harvey et al., 2024). However, it is difficult to understand the operations applied to the target data by the RBM to transform them into the source space due to their non-linear nature. Various network analysis methods could be used to test the sensitivity of the representation and decoding to specific inputs (Yang et al., 2022). So, where the goal of an analysis is prediction, classification, or other types of questions amenable to non-linear ML approaches, our RBM with Fisher divergence method provides a flexible, efficient framework for combining neural data across subjects, even when the recordings are sparse spiking representations.

4.4 Future work

Several extensions of the presented approach are possible and are currently part of our ongoing work. First, the derivation of the conditional distribution of the source features given the target features from the joint distribution represented by an RBM is a direction worthwhile pursuing as it relates to the cross-subject mapping problem and would avoid the issues related to noise-like initialization in the approach presented in this paper. Second, learning subject-invariant RBM models that can also predict feature representations of unseen target subjects is a key direction that should be pursued since it is directly related to the generalization capability of the approach. Finally, we are also investigating the application of the model to other neural signal modalities, including non-invasive modalities such as EEG signals.

Data availability statement

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and

accession number(s) can be found below: <https://duke.box.com/s/yxb8c59anm3m43dtxkyheuktlh8wp0vk>.

Ethics statement

The manuscript presents research on animals that do not require ethical approval for their study.

Author contributions

HY: Conceptualization, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing. MA: Conceptualization, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing. SW: Investigation, Software, Visualization, Writing – review & editing. JP: Data curation, Writing – review & editing. SS: Data curation, Funding acquisition, Resources, Supervision, Validation, Writing – review & editing. VT: Funding acquisition, Resources, Supervision, Validation, Writing – review & editing.

Funding

The author(s) declared that financial support was received for this work and/or its publication. This work was supported in part by Air Force Office of Scientific Research (AFOSR) Grant No. FA9550-22-1-0315. Marko Angelichinoski was supported in part by the Army Research Office Contract Number W911NF16-1-0368.

Conflict of interest

The author(s) declared that this work was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declared that generative AI was not used in the creation of this manuscript.

Any alternative text (alt text) provided alongside figures in this article has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of

their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fncom.2026.1710914/full#supplementary-material>

References

- Angjelichinoski, M., Choi, J., Banerjee, T., Pesaran, B., and Tarokh, V. (2020a). Cross-subject decoding of eye movement goals from local field potentials. *J. Neural Eng.* 17:016067. doi: 10.1088/1741-2552/ab6df3
- Angjelichinoski, M., Pesaran, B., and Tarokh, V. (2020b). Deep cross-subject mapping of neural activity. *arXiv [preprint]*. arXiv:2007.06407.
- Carreira-Perpinan, M. A., and Hinton, G. (2005). "On contrastive divergence learning," in *International workshop on Artificial Intelligence and Statistics (AISTATS)* (PMLR), 33–40.
- Chai, R., Ling, S. H., San, P. P., Naik, G. R., Nguyen, T. N., Tran, Y., et al. (2017). Improving EEG-based driver fatigue classification using sparse-deep belief networks. *Front. Neurosci.* 11:103. doi: 10.3389/fnins.2017.00103
- Chen, H., Jin, M., Li, Z., Fan, C., Li, J., He, H., et al. (2021). MS-MDA: multisource marginal distribution adaptation for cross-subject and cross-session EEG emotion recognition. *Front. Neurosci.* 15:778488. doi: 10.3389/fnins.2021.778488
- CTRL-labs at Reality Labs, Sussillo, D., Kaifosh, P., Reardon, T. (2024). A generic noninvasive neuromotor interface for human-computer interaction. *bioRxiv*. doi: 10.1101/2024.02.23.581779
- Dabagia, M., Kording, K. P., and Dyer, E. L. (2023). Aligning latent representations of neural activity. *Nat. Biomed. Eng.* 7, 337–343. doi: 10.1038/s41551-022-00962-7
- Degenhart, A. D., Bishop, W. E., Oby, E. R., Tyler-Kabara, E. C., Chase, S. M., Batista, A. P., et al. (2020). Stabilization of a brain-computer interface via the alignment of low-dimensional spaces of neural activity. *Nat. Biomed. Eng.* 4, 672–685. doi: 10.1038/s41551-020-0542-9
- Dyer, E. L., Gheshlaghi Azar, M., Perich, M. G., Fernandes, H. L., Naufel, S., Miller, L. E., et al. (2017). A cryptography-based approach for movement decoding. *Nat. Biomed. Eng.* 1, 967–976. doi: 10.1038/s41551-017-0169-7
- Fahimi, F., Zhang, Z., Goh, W. B., Lee, T.-S., Ang, K. K., Guan, C., et al. (2019). Inter-subject transfer learning with an end-to-end deep convolutional neural network for EEG-based BCI. *J. Neural Eng.* 16:026007. doi: 10.1088/1741-2552/aaf3f6
- Farahani, H. S., Fatehi, A., and Shoorehdeli, M. A. (2020). Between-domain instance transition via the process of GIBBS sampling in RBM. *arXiv [preprint]*. arXiv:2006.14538.
- Gallego, J. A., Perich, M. G., Chowdhury, R. H., Solla, S. A., and Miller, L. E. (2020). Long-term stability of cortical population dynamics underlying consistent behavior. *Nat. Neurosci.* 23, 260–270. doi: 10.1038/s41593-019-0555-4
- Hajinoroozi, M., Mao, Z., Jung, T.-P., Lin, C.-T., and Huang, Y. (2016). Eeg-based prediction of driver's cognitive performance by deep convolutional neural network. *Signal Process. Image Commun.* 47, 549–555. doi: 10.1016/j.image.2016.05.018
- Harvey, S. E., Larsen, B. W., and Williams, A. H. (2024). "Duality of bures and shape distances with implications for comparing neural representations," in *Proceedings of UniReps: the First Workshop on Unifying Representations in Neural Models* (PMLR), 11–26.
- Hasan, A., Yang, H., Ng, Y., and Tarokh, V. (2025). "Elliptic loss regularization," in *The Thirteenth International Conference on Learning Representations*.
- Haxby, J. V., Guntupalli, J. S., Connolly, A. C., Halchenko, Y. O., Conroy, B. R., Gobbini, M. I., et al. (2011). A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron* 72, 404–416. doi: 10.1016/j.neuron.2011.08.026
- Herrero-Vidal, P., Rinberg, D., and Savin, C. (2021). Across-animal odor decoding by probabilistic manifold alignment. *Adv. Neural Inf. Process. Syst.* 34, 20360–20372. doi: 10.1101/2021.06.06.447279
- Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Comput.* 14, 1771–1800. doi: 10.1162/089976602760128018
- Hjelm, R. D., Calhoun, V. D., Salakhutdinov, R., Allen, E. A., Adali, T., Plis, S. M., et al. (2014). Restricted boltzmann machines for neuroimaging: an application in identifying intrinsic networks. *NeuroImage* 96, 245–260. doi: 10.1016/j.neuroimage.2014.03.048
- Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *J. Mach. Learn. Res.* 6, 695–709. doi: 10.5555/1046920.1088696
- Hyvärinen, A. (2007). Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables. *IEEE Trans. Neural Netw.* 18, 1529–1531. doi: 10.1109/TNN.2007.895819
- Jayaram, V., Alamgir, M., Altun, Y., Scholkopf, B., and Grosse-Wentrup, M. (2016). Transfer learning in brain-computer interfaces. *IEEE Comput. Intell. Mag.* 11, 20–31. doi: 10.1109/MCI.2015.2501545
- Karpowicz, B. M., Ali, Y. H., Wimalasena, L. N., Sedler, A. R., Keshtkaran, M. R., Bodkin, K., et al. (2022). Stabilizing brain-computer interfaces through alignment of latent dynamics. *bioRxiv*. doi: 10.1101/2022.04.06.487388
- Kim, H.-C., Jang, H., and Lee, J.-H. (2020). Test-retest reliability of spatial patterns from resting-state functional MRI using the restricted Boltzmann machine and hierarchically organized spatial patterns from the deep belief network. *J. Neurosci. Methods* 330:108451. doi: 10.1016/j.jneumeth.2019.108451
- Kingma, D. P., and Ba, J. (2014). "ADAM: a method for stochastic optimization," in *The Third International Conference on Learning Representations*.
- Klabunde, M., Schumacher, T., Strohmaier, M., and Lemmerich, F. (2023). Similarity of neural network models: a survey of functional and representational measures. *arXiv [preprint]*. arXiv:2305.06329.
- Kriegeskorte, N., Mur, M., and Bandettini, P. A. (2008). Representational similarity analysis-connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* 2:249. doi: 10.3389/fnins.2008.004.2008
- Kuremoto, T., Kimura, S., Kobayashi, K., and Obayashi, M. (2014). Time series forecasting using a deep belief network with restricted Boltzmann machines. *Neurocomputing* 137, 47–56. doi: 10.1016/j.neucom.2013.03.047
- Lee, J., Dabagia, M., Dyer, E., and Rozell, C. (2019). "Hierarchical optimal transport for multimodal distribution alignment," in *Advances in Neural Information Processing System*, 32.
- Li, F., Tran, L., Thung, K.-H., Ji, S., Shen, D., Li, J., et al. (2015). A robust deep model for improved classification of ad/mci patients. *IEEE J. Biomed. Health Inform.* 19, 1610–1616. doi: 10.1109/JBHI.2015.2429556
- Li, X., Song, D., Zhang, P., Zhang, Y., Hou, Y., Hu, B., et al. (2018). Exploring EEG features in cross-subject emotion recognition. *Front. Neurosci.* 12:162. doi: 10.3389/fnins.2018.00162
- Monnens, S. Q., Peters, C., Hesselink, L. W., Smeets, K., and Englitz, B. (2024). The recurrent temporal restricted Boltzmann machine captures neural assembly dynamics in whole-brain activity. *eLife* 13:RP98489. doi: 10.7554/eLife.98489
- Pandarath, C., Ames, K. C., Russo, A. A., Farshchian, A., Miller, L. E., Dyer, E. L., et al. (2018). Latent factors and dynamics in motor cortex and their application to brain-machine interfaces. *J. Neurosci.* 38, 9390–9401. doi: 10.1523/JNEUROSCI.1669-18.2018
- Plis, S. M., Hjelm, D. R., Salakhutdinov, R., Allen, E. A., Bockholt, H. J., Long, J. D., et al. (2014). Deep learning for neuroimaging: a validation study. *Front. Neurosci.* 8:229. doi: 10.3389/fnins.2014.00229
- Putney, J., Angjelichinoski, M., Ravier, R., Ferrari, S., Tarokh, V., Sponberg, S., et al. (2021). Consistent coordination patterns provide near perfect behavior decoding in a comprehensive motor program for insect flight. *bioRxiv*. doi: 10.1101/2021.07.13.452211
- Putney, J., Conn, R., and Sponberg, S. (2019). Precise timing is ubiquitous, consistent, and coordinated across a comprehensive, spike-resolved flight motor program. *Proc. Natl. Acad. Sci.* 116, 26951–26960. doi: 10.1101/602961
- Rao, R. P. N. (2013). *Brain-Computer Interfacing: An Introduction*. Cambridge, MA: Cambridge University Press. doi: 10.1017/CBO9781139032803
- Shalit, U., Johansson, F. D., and Sontag, D. (2017). "Estimating individual treatment effect: generalization bounds and algorithms," in *International Conference on Machine Learning* (PMLR), 3076–3085.

- Sohn, K., Lee, H., and Yan, X. (2015). "Learning structured output representation using deep conditional generative models," in *Advances in Neural Information Processing Systems (NIPS)*, 28.
- Soleimani, E., and Nazerfard, E. (2021). Cross-subject transfer learning in human activity recognition systems using generative adversarial networks. *Neurocomputing* 426, 26–34. doi: 10.1016/j.neucom.2020.10.056
- Tian, H., and Xu, Q. (2021). Time series prediction method based on e-crbm. *Electronics* 10:416. doi: 10.3390/electronics10040416
- Tieleman, T. (2008). "Training restricted Boltzmann machines using approximations to the likelihood gradient," in *Proceedings of the 25th International Conference on Machine Learning*, 1064–1071. doi: 10.1145/1390156.1390290
- Torres-Oviedo, G., and Ting, L. H. (2010). Subject-specific muscle synergies in human balance control are consistent across different biomechanical contexts. *J. Neurophysiol.* 103, 3084–3098. doi: 10.1152/jn.00960.2009
- van der Plas, T. L., Tubiana, J., Le Goc, G., Migault, G., Kunst, M., Baier, H., et al. (2023). Neural assemblies uncovered by generative modeling explain whole-brain activity statistics and reflect structural connectivity. *Elife* 12:e83139. doi: 10.7554/eLife.83139
- Wang, Y., Wang, J., Wang, W., Su, J., Bunterngrchit, C., Hou, Z.-G., et al. (2024). TFFL: a task-free transfer learning strategy for EEG-based cross-subject cross-dataset motor imagery BCI. *IEEE Trans. Biomed. Eng.* 72, 810–821. doi: 10.1109/TBME.2024.3474049
- Wei, B., and Pal, C. (2011). Heterogeneous transfer learning with RBMs. *Proc. AAAI Confe. Artif. Intell.* 25, 531–536. doi: 10.1609/aaai.v25i1.7925
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica* 50, 1–25. doi: 10.2307/1912526
- Wu, Y., and Ji, Q. (2016). "Constrained deep transfer feature learning and its applications," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5101–5109. doi: 10.1109/CVPR.2016.551
- Yang, H., Winter, S., Zhang, Z., and Dunson, D. (2022). Interpretable ai for relating brain structural and functional connectomes. *arXiv [preprint]*. arXiv:2210.05672.
- Zhang, L., Xiao, D., Guo, X., Li, F., Liang, W., Zhou, B., et al. (2023). Cross-subject emotion EEG signal recognition based on source microstate analysis. *Front. Neurosci.* 17:1288580. doi: 10.3389/fnins.2023.1288580
- Zhao, Y., Dai, G., Borghini, G., Zhang, J., Li, X., Zhang, Z., et al. (2021). Label-based alignment multi-source domain adaptation for cross-subject EEG fatigue mental state evaluation. *Front. Hum. Neurosci.* 15:706270. doi: 10.3389/fnhum.2021.706270
- Zhong, Y., Yao, L., Pan, G., and Wang, Y. (2024). Cross-subject motor imagery decoding by transfer learning of tactile ERD. *IEEE Trans. Neural Syst. Rehab. Eng.* 32, 662–671. doi: 10.1109/TNSRE.2024.3358491