
Learning identifiable and interpretable latent models of high-dimensional neural activity using pi-VAE

Ding Zhou

Department of Statistics
Columbia University
dz2336@columbia.edu

Xue-Xin Wei

Department of Neuroscience
UT Austin
weixx@utexas.edu

Abstract

The ability to record activities from hundreds of neurons simultaneously in the brain has placed an increasing demand for developing appropriate statistical techniques to analyze such data. Recently, deep generative models have been proposed to fit neural population responses. While these methods are flexible and expressive, the downside is that they can be difficult to interpret and identify. To address this problem, we propose a method that integrates key ingredients from latent models and traditional neural encoding models. Our method, pi-VAE, is inspired by recent progress on identifiable variational auto-encoder, which we adapt to make appropriate for neuroscience applications. Specifically, we propose to construct latent variable models of neural activity *while simultaneously* modeling the relation between the latent and task variables (non-neural variables, *e.g.* sensory, motor, and other externally observable states). The incorporation of task variables results in models that are not only more constrained, but also show qualitative improvements in interpretability and identifiability. We validate pi-VAE using synthetic data, and apply it to analyze neurophysiological datasets from rat hippocampus and macaque motor cortex. We demonstrate that pi-VAE not only fits the data better, but also provides unexpected novel insights into the structure of the neural codes.

1 Introduction

Popular analysis methods of neural responses in neurophysiology mainly come in two classes: one based on regression, the other on latent variable modeling. Generalized Linear Model (*e.g.*, [67, 56]) and tuning curve analysis [24, 73, 50] are notable examples of the regression-based approach, and both have been widely used in neuroscience in the past few decades [63, 60, 73, 54, 11]. These methods express the neural firing rate as a function of the stimulus variable, thus naturally define encoding models, which can be inverted to decode the stimulus variables [24, 59, 73, 6, 51]. In contrast, the latent-based approach aims to account to variability of the neural responses using a relatively small number of latent variables which are typically not observed. Recently, various latent-based methods have been developed or applied to analyze neural data and in particular simultaneously recorded neural population data, including principal component analysis [66, 46, 4, 13, 45], factor analysis [10, 58, 18], linear/nonlinear dynamical systems [44, 7, 22, 17, 52], among others (*e.g.*, [5, 74, 71, 34]). Each class of models carries certain advantages and disadvantages. Regression-based methods tend to have higher interpretability, however, they often suffer the problem of under-fitting. Latent-based models are more flexible in accounting for the neural variability, however they may be difficult to interpret and sometimes not identifiable. Notably, some studies had incorporated latent fluctuations into the encoding models [72, 39, 25, 19, 42, 11] yielding promising results, although these models often assumed highly specialized latent structure, thus potentially limits the applicability in practice.

The issues of identifiability and interpretability are becoming increasingly important as the neuroscience community adapts more sophisticated methods from nonlinear deep generative models [17, 70]. Deep generative models have the promise of extracting complex nonlinear structure which may be difficult to achieve by linear methods, as demonstrated by recent work based on the variational auto-encoder (VAE) (e.g., [37, 57, 22, 52]). However, over the past few years it has become increasingly clear that the latents extracted from these models, and VAE in particular, are often highly entangled therefore difficult to interpret [29, 28, 1, 12, 43, 8, 35]. Given these considerations, an important question is how to model neural population responses with nonlinear models that are powerful yet scientifically insightful via identifiability and interpretability.

We propose a model formulation which represents one step toward addressing this question. Specifically, we draw on recent progress on identifiable VAE (iVAE) [35, 64], and generalize and adapt it to make it directly applicable to a broad variety of datasets. Conceptually, our method combines the respective strengths of regression-based and latent-based approaches (Fig. 1): i) by using the VAE architecture, it is expressive and flexible; ii) by treating task variables as labels and explicitly modeling them with the latents, our method is better constrained and under certain conditions identifiable. We apply our method to synthetic data and electrophysiological datasets from population recording of rat hippocampus in a navigation task [27, 26] and the motor area of macaque during a reaching task [21]. We demonstrate that our method can recover interpretable latent structure that is informative about the structure of the neural code and dynamics.

2 Model

Notations We denote $\mathbf{x} \in \mathbb{R}^n$ as observations. For our purpose, \mathbf{x} specifically represents the population response or spike counts within a small time window. We use $\mathbf{u} \in \mathbb{R}^d$ to represent *task variables* (or *labels*), that are measured along with the neural activities, e.g., the location of the animal when studying navigation tasks. \mathbf{u} can be discrete or continuous. Additionally, we denote $\mathbf{z} \in \mathbb{R}^m (m \ll n)$ as the unobserved low-dimensional latent variables.

2.1 Generative model

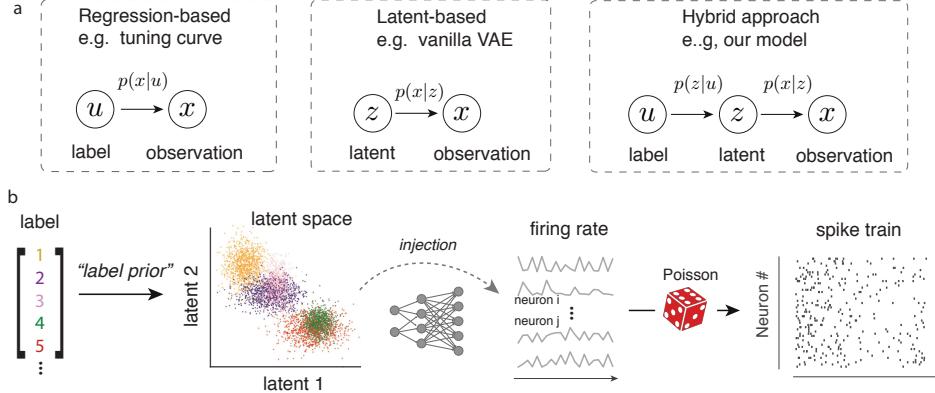


Figure 1: The model framework and generative model. (a) Structure of three classes of statistical models for neural data analysis. Our method is based on the integration of the first two classes into a hybrid approach. Our approach models the statistical dependence between label (\mathbf{u}) and latent (\mathbf{z}) as well as between latent (\mathbf{z}) and observation (\mathbf{x}) simultaneously. (b) Schematic illustration of the generative model of pi-VAE. Major components include the “label prior” between the task variables and the latent, an injective mapping between latent and firing rate parameterized by normalizing flow, and Poisson observation noise.

Our goal is to develop models that are flexible and expressive in capturing the variability of the data, while also well-constrained so that the models would enjoy identifiability and interpretability. Motivated by these considerations, we propose a generative model formulation which integrates key ingredients of the latent-based and regression-based approaches (see Fig. 1a):

$$p_{\theta}(\mathbf{x}, \mathbf{z}|\mathbf{u}) = p_f(\mathbf{x}|\mathbf{z})p_{T,\lambda}(\mathbf{z}|\mathbf{u}). \quad (1)$$

This is a general formulation, and some previous models may be re-formulated to conform with it (*e.g.*, [39, 25, 19]). In this paper, we will focus on a specific implementation that is directly inspired by the recent work on identifiable VAE [35, 41, 64]. In the interest of neuroscience applications, we have developed a method that can simultaneously deal with Poisson noise, both discrete and continuous labels, and **larger output dimension than input dimension** (Fig. 1b). We will show that our model is sufficiently expressive for many applications, yet still constrained enough to be identifiable.

We start by defining the component that describes the relation between the label and the latent, *i.e.*, $p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{z}|\mathbf{u})$. We will refer to it as the “*label prior*”. Following [35], we assume $p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{z}|\mathbf{u})$ to be conditionally independent, where each element $\mathbf{z}_i \in \mathbf{z}$ has an *exponential family distribution* given \mathbf{u} ,

$$p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{z}|\mathbf{u}) = \prod_{i=1}^m p(\mathbf{z}_i|\mathbf{u}) = \prod_{i=1}^m \frac{Q_i(\mathbf{z}_i)}{Z_i(\mathbf{u})} \exp \left[\sum_{j=1}^k T_{i,j}(\mathbf{z}_i) \lambda_{i,j}(\mathbf{u}) \right], \quad (2)$$

where Q_i is the base measure, $\mathbf{T}_i = (T_{i,1}, \dots, T_{i,k})$ are the sufficient statistics, $Z_i(\mathbf{u})$ is the normalizing factor, $\boldsymbol{\lambda}_i = (\lambda_{i,1}(\mathbf{u}), \dots, \lambda_{i,k}(\mathbf{u}))$ are the natural parameters, and k is pre-defined number of sufficient statistics. Practically, a small number of components k is often sufficient for the problems which we have considered. For discrete \mathbf{u} , we simply use a different λ_{ij} for different \mathbf{u} . To deal with continuous \mathbf{u} , we develop a procedure by parameterizing λ_{ij} as a function of \mathbf{u} using a feed-forward neural network. Details are given in Supplementary Information (SI).

We next turn to the dependence between the latent and observation, *i.e.*, $p_{\mathbf{f}}(\mathbf{x}|\mathbf{z})$. In [35], $p_{\mathbf{f}}(\mathbf{x}|\mathbf{z})$ is defined using additive noise $p_{\mathbf{f}}(\mathbf{x}|\mathbf{z}) = p_{\epsilon}(\mathbf{x} - \mathbf{f}(\mathbf{z}))$, *i.e.* $\mathbf{x} = \mathbf{f}(\mathbf{z}) + \epsilon$, where ϵ is an independent noise variable. To model the spike data, we generalize it to Poisson model $p_{\mathbf{f}}(\mathbf{x}|\mathbf{z}) = \text{Poisson}(\mathbf{f}(\mathbf{z}))$ with \mathbf{f} being the instantaneous firing rate to deal with the count observations. We implement \mathbf{f} using normalizing flow as detailed later. Putting together, we denote $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$ as parameters in generative model 1. We refer to our model as Poission identifiable VAE, or pi-VAE for simplicity.

Identifiability [35] has proved that the additive noise model is identifiable with certain assumptions. Under the same assumptions, we can prove that pi-VAE is also identifiable.

Definition 1. Let \sim be an equivalence relation on the domain of the parameters $\Theta = (\boldsymbol{\theta} := (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}))$. Model 1 is said to be identifiable up to \sim if $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{u}) = p_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}|\mathbf{u}) \implies \boldsymbol{\theta} \sim \tilde{\boldsymbol{\theta}}$.

Definition 2. Define \sim as $\boldsymbol{\theta} \sim \tilde{\boldsymbol{\theta}} \iff \exists A, c, \mathbf{T}(\mathbf{f}^{-1}(\mathbf{x})) = A\tilde{\mathbf{T}}(\tilde{\mathbf{f}}^{-1}(\mathbf{x})) + c, \forall \mathbf{x} \in \text{Img}(\mathbf{f}) \subset \mathbb{R}^n$, where $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}), \tilde{\boldsymbol{\theta}} = (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}})$, A is a full rank $mk \times mk$ matrix, $c \in \mathbb{R}^{mk}$ is a vector.

Theorem 1. Assume that we observe data sampled from pi-VAE model defined according to equation 1, 2 with Poisson noise and parameters $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$. Assume the following holds:

- i) The firing rate function \mathbf{f} in equation 1 is injective,
- ii) The sufficient statistics $T_{i,j}$ in 2 are differentiable almost everywhere, and their derivatives $T'_{i,j}$ are nonzero almost everywhere for $1 \leq i \leq m, 1 \leq j \leq k$,
- iii) There exists $mk + 1$ distinct points $\mathbf{u}^0, \dots, \mathbf{u}^{mk+1}$ such that the matrix

$$L = (\boldsymbol{\lambda}(\mathbf{u}^1) - \boldsymbol{\lambda}(\mathbf{u}^0), \dots, \boldsymbol{\lambda}(\mathbf{u}^{mk}) - \boldsymbol{\lambda}(\mathbf{u}^0))$$

of size $mk \times mk$ is invertible, then the pi-VAE model is identifiable up to \sim .

This theorem is a straight-forward generalization of the results in [35]. Proof is given in the SI. This theorem means, if two sets of model parameters lead to the same marginal distribution of \mathbf{x} , then $(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}) \sim (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}})$. Thus one can hope to recover posterior distribution $p(\mathbf{z}|\mathbf{x})$ up to a linear transformation A and point-wise non-linearities between \mathbf{T} and $\tilde{\mathbf{T}}$, as well as the joint distribution $p(\mathbf{x}, \mathbf{z})$. Note that other forms of identifiability may be derived with modified assumptions [35]. While practically, without knowing the ground truth, the assumptions may be difficult to verify, some encouraging preliminary evidence suggest that identifiability may have some robustness with mild violations of model assumptions [64].

We next describe how to parameterize the injection $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$. Here we extend the General Incompressible-flow Network (GIN) proposed in [64], which shares the flexibility of RealNVP [16] and the volume-preserving property of NICE [15]. Practically, we found our implementation to be reasonably efficient computationally. Specifically, GIN defines a mapping from $\mathbb{R}^D \rightarrow \mathbb{R}^D$ with

Jacobian determinant equal to 1 [64]. It splits the D -dimensional input \mathbf{x} into two parts $\mathbf{x}_{1:l}, \mathbf{x}_{l+1:D}$, where $l < D$. The output \mathbf{y} is defined as the concatenation of $\mathbf{y}_{1:l}$ and $\mathbf{y}_{l+1:D}$,

$$\mathbf{y}_{1:l} = \mathbf{x}_{1:l} \quad (3)$$

$$\mathbf{y}_{l+1:D} = \mathbf{x}_{l+1:D} \odot \exp(s(\mathbf{x}_{1:l})) + t(\mathbf{x}_{1:l}), \quad (4)$$

where $s(\cdot)$ and $t(\cdot)$ are both functions defined on $\mathbb{R}^l \rightarrow \mathbb{R}^{D-l}$, and the total sum of $s(\mathbf{x}_{1:l})$ is constrained to be zero by setting the final component to be the negative sum of previous components.

The original GIN [64] only deals with the case where the input and output have the same dimensions. In our case, the output dimension is often much larger than the input dimension. We thus develop a new scheme to parameterize $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ which retains the properties of GIN. We first map $\mathbf{z}_{1:m}$ to the concatenation of $\mathbf{z}_{1:m}$ and $t(\mathbf{z}_{1:m})$. Note that this is equivalent to GIN model with input as $\mathbf{z}_{1:m}$ padding $n-m$ zeros. We then use several GIN blocks to map from $\mathbb{R}^n \rightarrow \mathbb{R}^n$. Recalling that each GIN block is an injection (since part of it is an identity map, one can not map two different inputs to the same output), it follows that the composition of several blocks remains an injection. While we use an extension of GIN [64] to implement the injection f here, conceivably other implementations should be possible, e.g., multi-layer perceptron with increasing number of nodes from earlier to later layers. The efficiency of the different implementations will need to be evaluated in future.

2.2 Inference algorithm

The inference procedure is a modification of VAE [37]. Our algorithm simultaneously learns the deep generative model and the approximate posterior $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$ of true posterior $p(\mathbf{z}|\mathbf{x}, \mathbf{u})$ by maximizing $\mathcal{L}(\theta, \phi)$, which is the evidence lower bound (ELBO) of $p(\mathbf{x}|\mathbf{u})$,

$$\log p(\mathbf{x}|\mathbf{u}) \geq \mathcal{L}(\theta, \phi) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{u})} [\log(p(\mathbf{x}, \mathbf{z}|\mathbf{u})) - \log(q(\mathbf{z}|\mathbf{x}, \mathbf{u}))]. \quad (5)$$

Similar to [31], we decompose the approximate posterior as

$$q(\mathbf{z}|\mathbf{x}, \mathbf{u}) \propto q_\phi(\mathbf{z}|\mathbf{x}) p_{\mathbf{T}, \lambda}(\mathbf{z}|\mathbf{u}), \quad (6)$$

where $q_\phi(\mathbf{z}|\mathbf{x})$ is assumed to be conditionally independent exponential family distribution, i.e. $q_\phi(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^m q(\mathbf{z}_i|\mathbf{x})$, and is parameterized by neural network. $p_{\mathbf{T}, \lambda}(\mathbf{z}|\mathbf{u})$ is defined in equation 2.

We modeled both $q_\phi(\mathbf{z}|\mathbf{x})$, $p_{\mathbf{T}, \lambda}(\mathbf{z}|\mathbf{u})$ as independent Gaussian distribution, used the same network architecture (see SI for details) as well as Adam optimizer [36] with learning rate equal to 5×10^{-4} , and other values were set to the recommendation values for all the experiments in this paper.

Inferred the latent After learning $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$, the latent from pi-VAE model can be inferred by calculating the posterior mean. It is also of interest to infer the latent without using the label prior $p_{\mathbf{T}, \lambda}(\mathbf{z}|\mathbf{u})$, which could be done by calculating the posterior mean estimate of $q_\phi(\mathbf{z}|\mathbf{x})$ instead.

Decoding the label Because pi-VAE defines an encoding model on the label, one can examine how well the label could be decoded from the neural activity, which also provides a way to check the validity of the model. Under our model formulation, decoding could be done by Bayesian rule and Monte Carlo sampling: $p(\mathbf{u}|\mathbf{x}) \propto \int p(\mathbf{x}|\mathbf{z}, \mathbf{u}) p(\mathbf{z}|\mathbf{u}) p(\mathbf{u}) d\mathbf{z}$, where the integration on right hand side can be computed through randomly sampling $p(\mathbf{z}|\mathbf{u})$. We assume a uniform prior on \mathbf{u} .

3 Results

3.1 Validation using synthetic data

We validated pi-VAE using synthetic data generated from models with continuous or discrete labels.

Discrete label We generated 2-dimensional latent samples \mathbf{z} from a five clusters Gaussian mixture model, similar to [35] (see Fig. 2a). The mean of each cluster was chosen independently from a uniform distribution on $[-5, 5]$ and variance from a uniform distribution on $[0.5, 3]$. These latent samples were then mapped to the mean firing rate of 100-dimensional Poisson observations through a RealNVP network (details in SI). Example results are shown in Fig. 2a-d.

Continuous label We generated \mathbf{u} from a uniform distribution on $[0, 2\pi]$, and latent samples \mathbf{z} as a 2-dimensional independent Gaussian distribution with mean being $(\mathbf{u}, 2 \sin \mathbf{u})$, and variance being

$(0.6 - 0.3|\sin \mathbf{u}|, 0.3|\sin \mathbf{u}|)$. Observations were generated in the same way as simulation of discrete label. Example results are shown in Fig. 2e-h.

Based on these and other numerical experiments, we found that in general pi-VAE could reliably uncover latent structure similar to ground truth for both discrete and continuous labels, while VAE often leads to more distorted latent. Note that our VAE implementation is similar to pi-VAE except that no label prior is used (*e.g.*, Poisson observation noise is still assumed). We also found pi-VAE without label prior (during the inference) still led to reasonably good recovery of the latent (Fig. 2c,g), suggesting that incorporating label prior could help with learning a better model, not just inference.

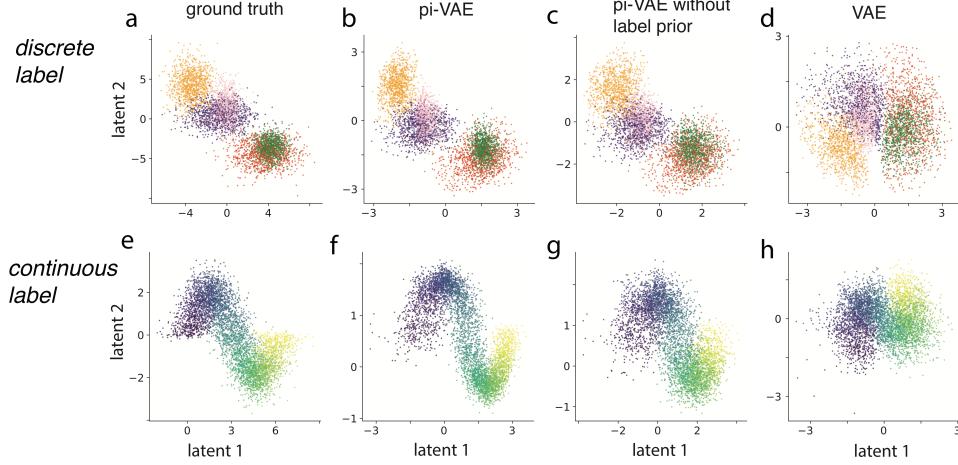


Figure 2: Example numerical experiments, suggesting that pi-VAE, but not VAE, could identify latent structure. (a) True latent variables, simulated based on discrete label, (b) mean of the latent posterior $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$ estimated from pi-VAE, (c) mean of $q(\mathbf{z}|\mathbf{x})$ from pi-VAE, (d) mean of the latent posterior from VAE, that is, the Bayesian estimate *inferred* without the label prior. (e-h) similar to (a-d), but for a simulation based on continuous label.

3.2 Applications to neural population data

We have applied pi-VAE to analyze two electrophysiology datasets, each has more than 100 simultaneously recorded neurons when the animals were performing behavioral tasks. In these real data applications the ground truth is unknown and the assumptions required by identifiability may be violated [64], making it difficult to assess identifiability directly. Our rationale is that, **assuming the ground truth is structured, models with better identifiability would still lead to more interpretable latent representation**. Encouragingly, examination of the latent space extracted from these datasets indeed suggest that pi-VAE could extract interpretable and meaningful latent structure.

3.2.1 Monkey reaching data

We first applied our method to a previously published monkey reaching datasets (see [21], kindly shared by the authors). In these experiments, Monkey C was performing a reaching task with 8 different directions, while neural activities in areas M1 and PMd were simultaneously recorded (for details, see [21]). We analyzed two sessions, and obtained similar results. We will focus on Session 1 here, and detailed results from Session 2 can be found in SI.

For each direction, there are ~ 25 trials/repeats (see Fig. 3a). We analyzed 192 neurons from PMd area, and focused on the reaching period from go cue (defined as $t = 0$) to the end, which typically lasts for ~ 1 second. We binned the ensemble spike activities into 50ms bins. We used the spike activities as observation \mathbf{x} , and the reaching direction as the discrete labels \mathbf{u} . We randomly split the dataset into 24 batches, where each batch contains at least one trial for each direction. We randomly split them into training, validation and test data (20, 2, 2 batches). We fit 4-dimensional latent models to the data based on pi-VAE and VAE.

Goodness of fit We first assessed the goodness of fit by examining the root-mean-square error (RMSE) of the PSTH based on the prediction of each model (see Fig. 3a for an example neuron).

Fig. 3b,c show that pi-VAE leads to the smallest RMSE of firing rate in most neurons, followed by VAE, then tuning curve model. Next, we computed the log marginal likelihood $p(\mathbf{x})$ on the held-out test data by randomly sampling both $p(\mathbf{z}|\mathbf{u})$ and $p(\mathbf{u})$. We found that pi-VAE leads to larger mean marginal log-likelihood than VAE and tuning curve model (-123 , -123.4 , -127.6 respectively, t-test $p < 10^{-6}$). These results suggest that pi-VAE provides the best fit to the data among the alternatives.

Decoding reaching direction We wondered whether pi-VAE also provides a better encoding model of the reaching direction. We examined how well pi-VAE could decode reaching direction, and compare the performance to a traditional method based on direction tuning curves. On held-out test data, pi-VAE achieved an average single-frame (50ms) decoding accuracy of 61%, better than 47% from tuning curve model. Examination of the time course of the decoding performance (Fig. 3d) reveals that pi-VAE achieves 60% during the first few frames before initiation of reaching, while tuning curve model is much worse during this period. However, when reaching speed reaches its maximum (around 0.5s, Fig. 3e), both models achieve almost perfect performance (Fig. 3d). These observations tentatively suggest that information about the reaching direction may be encoded in different format during different phases of reaching in this task.

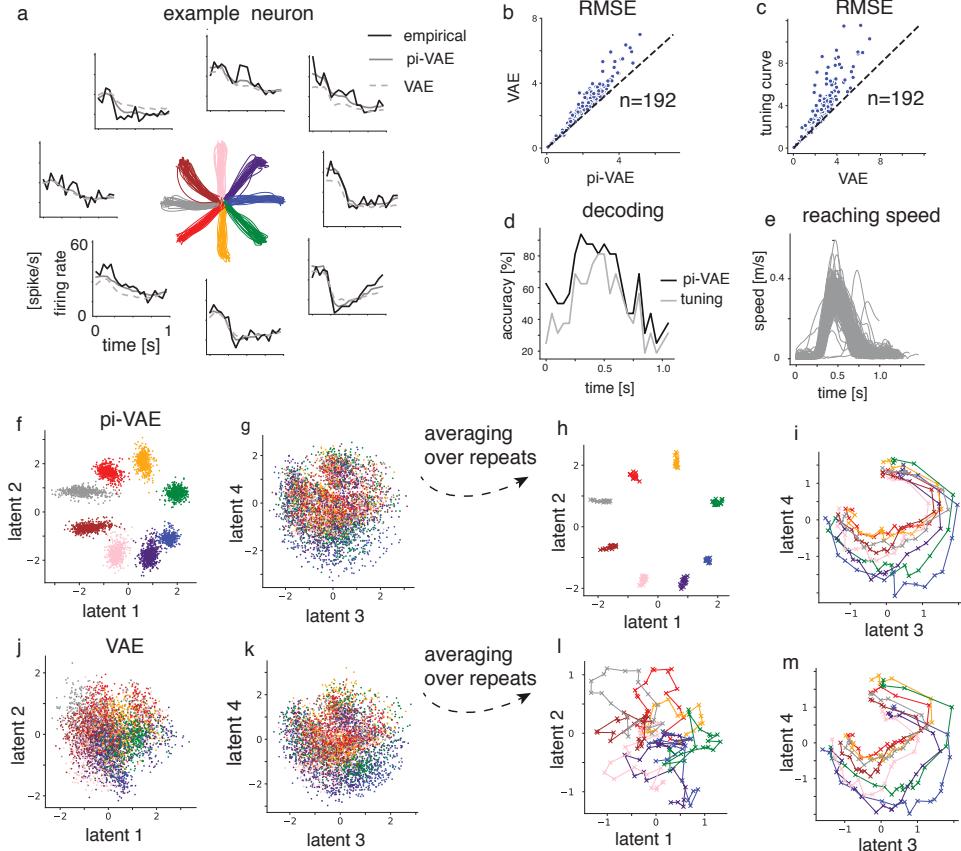


Figure 3: Monkey reaching data. (a) Reaching trajectories for 8 directions labeled by colors. The empirical firing rate (PSTH, black solid line), fitted rate by pi-VAE (gray solid line) and VAE (gray dashed line) for an example neuron. (b,c) Scatter plots of RMSE of fitted rate ($n = 192$ neurons) for comparing pi-VAE and VAE, as well as VAE and tuning curve. (d) Decoding accuracy as function of time on test data by pi-VAE and tuning curve model. (e) The reaching speed of the macaque for each trial. (f,g) Inferred latent based on pi-VAE, i.e.,mean of $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$. (h,i) Inferred latent from pi-VAE averaged over repeats from the same reaching direction. (j,k,l,m) Similar to (f,g,h,i) for VAE. Notice the striking difference between (f) and (j).

Structure of the latent We found that the latent variables estimated by pi-VAE exhibit clear structures. To start, the 8 reaching directions are well separated in the subspace defined by the first two latent dimensions (Fig. 3f,h). Strikingly, the geometrical structure of the inferred latent resembles the

geometry of the reaching directions. In contrast, the third and fourth latent dimensions captures the evolution of the trajectories over time, and they are only weakly informative about reaching directions (Fig. 3g,i). Thus, the axes of the extracted latent space are easily interpretable. They provide information about how reaching direction is represented, and how neural dynamics evolve during reaching behavior. Notably, these axes were extracted automatically from pi-VAE, and no additional factor analysis techniques were applied to identify salient latent axes. Strikingly, the latent variables extracted from Session 2 show very similar structure (see SI Fig. S1).

In comparison, VAE extracts a much more entangled latent representation (Fig. 3j-m). It appears that information about reaching direction displays in a twisted fashion, and mixes with the temporal evolution of the trajectories. Note that these differences between the latent structure obtained from two methods is not simply due to the label prior. The inferred latent from pi-VAE without the label prior (*i.e.*, posterior mean of $q(\mathbf{z}|\mathbf{x})$) shows similar though a bit more diffuse latent structure, which is expected due to the observation noise (see SI Fig. S2).

The nature of the neural code during primate reaching behavior is currently under heavy debate [24, 23, 60, 65, 13, 61, 40, 21]. While earlier proposals emphasized the encoding of task relevant variables [24, 23, 60, 53], some of the more recent studies instead highlighted the importance of neural dynamics [13, 61, 2, 32]. As shown above, pi-VAE discovers latent space that exhibits striking spatial (*i.e.*, reaching direction) *and* temporal (*i.e.*, neural dynamics) structure that are separately encoded in different sub-spaces. These preliminary results may provide a way to reconcile the two prominent hypotheses [24, 61], as evidence for both hypotheses are now revealed in the same model based on the same datasets. It would be important to apply our methods to larger datasets from multiple monkeys to examine the consistency of these effects in future.

3.2.2 Rat hippocampal CA1 data

We next applied pi-VAE to analyze a public rat's hippocampus dataset [27, 26]¹. In this experiment, a rat ran on a 1.6m linear track with rewards at both ends (L&R) (Fig. 4a), while neural activity in the hippocampal CA1 area was recorded ($n = 120$, putative pyramidal neurons). We focused on the data when the rat was running on the track (Fig. 4a) and binned the ensemble spike activities into 25ms bins. We defined the rat running from one end of the track to the other end as one lap, resulting in 84 laps. We randomly split them into training, validation and test data (68, 8, 8 laps). We defined rat's position and running directions as continuous labels \mathbf{u} . We fit 2-dimensional latent models to the data for both pi-VAE and VAE.

Goodness of fit and decoding performance We found that pi-VAE again outperformed alternatives in having the lowest mean log marginal likelihood -17.7 (VAE, -17.9 ; tuning curve, -18.2 ; paired t-test, $p < 10^{-6}$). Furthermore, we decoded the animal's location on the tracking based on pi-VAE model and tuning curve model. On test data, pi-VAE achieves median absolute decoding error (MAE) of 12cm (time window = 25ms), while the tuning curve (traditional "place field" [49]) model achieves a MAE of 15cm. This indicates that for the simple purpose of constructing an effective encoding model of the animal's position on the track, pi-VAE outperforms the traditional place field model [73].

Structure of the latent Fig. 4b shows the latent space estimated by the pi-VAE, which exhibits overtly interpretable geometry: the collection of inferred latent states for R-to-L (blue) or L-to-R (red) running direction each forms band-like sub-manifold, and both are roughly in parallel with the second latent dimension (Fig. 4b). The split into two sub-manifolds is consistent with the observation that place fields of CA1 neurons often have firing fields that are uncorrelated between the two travel directions ("directional" firing) [3, 14, 48]. To further quantify the geometrical relation between two sub-manifolds, we calculated the distance for pairs of points from the two branches (Fig. 4b). This quantification for every possible pair (after binning the position into 16cm bins) is plotted in Fig. 4c. We found that the manifold geometry across the two directions respects the geometry of the track, in the sense that smaller physical distance on the track leads to smaller latent distance. This is likely due to that a subset of place cells have non-directional (purely spatial) place fields [3, 14, 48, 33]. Importantly, our method gives a quantitative population level characterization of the consequence of having both directional and non-directional place cells. pi-VAE without label prior shows similar but more diffuse latent structure (see SI Fig. S3). In contrast, VAE results in a tangled latent representation, with the geometry not reflecting physical distance on the track (Fig. 4d,e; also notice that striking difference between Fig. 4b and Fig. 4d).

¹<http://crcns.org/data-sets/hc/hc-11>

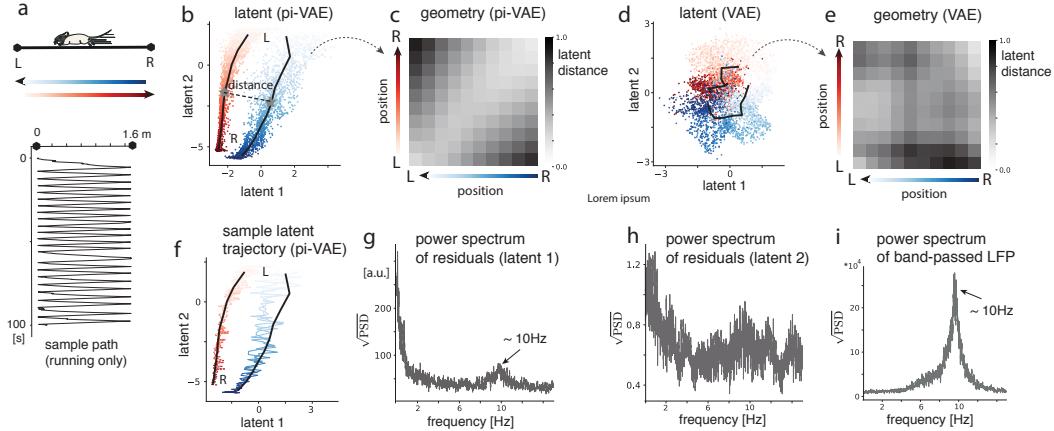


Figure 4: Results on hippocampus CA1 data. (a) Linear track and sample running path. The two ends are labeled as L&R. Two directions are color-coded by red and blue, and positions are coded by color saturation. (b) Inferred latent from pi-VAE. Black lines represent the mean of the latent states corresponding to position on the track for two directions. The distance between pairs of points from the two black lines is computed to quantify the latent geometry. An example pair of points are indicated using grey stars. The normalized distance for all possible pairs of points is shown in panel (c). (d,e) are defined similar to (b,c) for the VAE. Notice the striking difference between (b) and (d). (f) Two sample latent trajectories of the pi-VAE. (g,h) The power spectrum (PSD) of the residuals of the latent, given by the mean of $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$ minus the mean of $p(\mathbf{z}|\mathbf{u})$. (i) PSD of the band-passed LFP in the range of 5-11 Hz.

To investigate whether the latent model could yield additional scientific insight, we next examined the temporal structure of the latent. Fig. 4f plots sample trajectories, from which we observed that temporal fluctuation mainly goes along the first latent dimension, and the suggestion of rhythmic structure. We subtracted the mean of prior $p(\mathbf{z}|\mathbf{u})$ from the mean of posterior $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$ to obtain the residual fluctuations. Examination of the power spectrum density (PSD) along each dimension of the residuals led to two observations: i) the temporal fluctuation is indeed concentrated on the first latent dimension, as indicated by the magnitude of the PSD; ii) the first, but not the second, dimension exhibits a striking peak at ~ 10 Hz. We reasoned that the second observation might be related to θ -oscillation in the local circuit, which is known to modulate the firing of CA1 neurons [62, 20, 55, 33]. We thus examined the simultaneously recorded local field potential (LFP) data during running. Indeed, we found that the θ peaked at ~ 10 Hz for this rat. Interestingly, the 10Hz θ -oscillation is faster than the typically reported 8 Hz average frequency [69, 9], yet is consistent with the latent structure extracted from pi-VAE.

Overall, pi-VAE extracts latent space that is clearly interpretable, with one dimension encoding position information, and the other dimension capturing temporal organization which is likely related to θ -rhythm. The observation that the position encoding and the rhythmic-like fluctuation are roughly orthogonal is particularly interesting, and is consistent with previous results from the single cell analysis suggesting that theta phase and place fields may encode independent information [30]. Additional investigations will be needed to test these hypotheses in greater depth.

3.3 Comparison to alternative methods

We further tested several alternative methods on the monkey reaching data, including both linear methods (demixed PCA [38]) and nonlinear methods (UMAP [47], PfLDS [22], and LFADS [52]) (see Fig. 5). Overall we found that, while the extracted latent structures from these methods exhibited interesting characteristics, none of them resulted in fully disentangled latents. Furthermore, none of them appeared to recover the geometry of the physical reaching targets. (More analysis of the hippocampus data can be found in SI Sec. E)

To start, supervised UMAP recovered latents corresponding to different directions as different clusters, but without clear representations of temporal dynamics (see Fig. 5a). Furthermore, LFADS [52] and

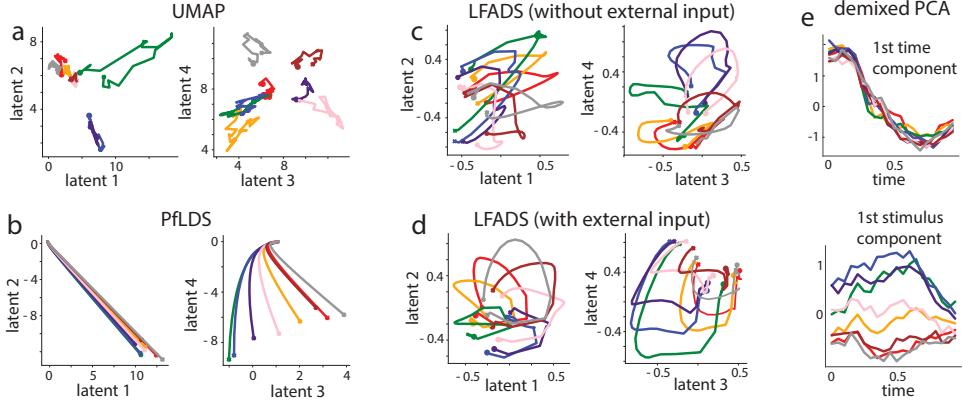


Figure 5: Results from alternative methods based on monkey reaching data. a) UMAP [47]. b) PfLDS [22]. c,d) LFADS [52], without and with reaching direction as an external input. e) demixed PCA [38] with the first time and stimulus component plotted. Color-coded, averaged latent trajectories corresponding to each reaching direction was plotted for each method. The filled dot and cross represent starting and ending of the trial.

PfLDS [22] both led to smooth trajectories. Although the trajectories for different directions were separated in the 4-dimensional space, directions and temporal dynamics were entangled so that it was difficult to interpret each individual latent dimension (Fig. 5b,c,d). Demixed PCA [38] with both time and directions as labels still entangled time and directions (stimulus components change with time) to some extent (Fig. 5e). A few methodological considerations are worth mentioning here. First, LFADS can take task variables as external inputs to the model RNN. We thus tried LFADS with or without reaching direction as external inputs (Fig. 5c,d). Second, demixed PCA only deals with discrete task variables each with the same number of trials and each trial with the same length, and could not recover additional latent fluctuations as our method. Third, UMAP can incorporate label information for supervised learning, and we used the reaching directions as labels (Fig. 5c,d) to make a more fair comparison. However, we found that it did not recover temporal dynamics.

4 Discussion

We have presented a new model framework for analyzing neural population data by integrating ingredients from latent-based and regression-based approaches. Our model pi-VAE, while being expressive and nonlinear, is constrained by additional dependence on task variables. pi-VAE generalizes recent work on identifiable VAE [35, 41, 64] to deal with spike train data. Although pi-VAE yields promising preliminary insights into the neural codes during a rat navigation task and macaque reaching task, we should emphasize that more systematic investigations based on larger datasets across different subjects will be needed to further elaborate these results.

Our method is motivated by leveraging the strength of regression-based methods and latent-based models to increase the identifiability and interpretability, a direction received little attention previously. To do so, we took advantage of the “label prior” to model the impact of task variables on neural activities along with the influence of the latent states. One potential concern is that, when incorporating too many labels, there may not be enough data to fit the model. Several previous methods exploited temporal smoothness priors to de-noise the data, which were implemented via Gaussian process [10, 74, 71], linear [44, 7, 22] or nonlinear dynamical systems [52, 17]. Although not pursued here, adding temporal smoothness priors into pi-VAE may increase the data efficiency and further improve the performance of the model. It is also worth mentioning that although we have focused on the spike train data, our method may be modified to deal with the calcium imaging data incorporating noise models that is more appropriate to the deconvolved calcium traces [68]. Last but not least, while the current study mainly concerns the neuroscience applications of pi-VAE, some of the technical advances made here may be of interest to the machine learning community as well.

Acknowledgements We thank Matthew G. Perich and Lee E. Miller for sharing the monkey reaching data. We thank Kenneth Kay, Yuanjun Gao and Liam Paninski for helpful comments and discussions, as well as three anonymous reviewers for their feedback which helped improving the paper. Xue-Xin Wei is supported by the startup funds provided by UT Austin.

5 Broader Impact

In this paper, we develop a new analysis method for analyzing neural population data. This method can be used to extract and identify the underlying critical structure underlying the neural activity recorded simultaneously from many neurons in the brain. It can also be used to construct improved response models of how various kinds of task variables are encoded in the brain. Our approach is of broad applicability to the neural population recording under a variety of scenarios. It may also inspire future work in neural data analysis. Progress in this area will facilitate both basic science research about the brain as well as clinical applications, such as invasive and non-invasive brain-machine interfaces.

Our method is developed based on recent advances in a class of deep generative models, i.e., identifiable variational auto-encoder. It may be of interest to the machine learning community working on deep generative models. Our research extends previous works on identifiable variational auto-encoder to a different setup. Our preliminary promising results may provide useful insights into the further development of more identifiable deep generative models, which in the longer term may help establish more interpretable and robust AI technology, and help understand the limitations of these technologies.

References

- [1] Alexander A Alemi, Ben Poole, Ian Fischer, Joshua V Dillon, Rif A Saurous, and Kevin Murphy. Fixing a broken elbo. *arXiv preprint arXiv:1711.00464*, 2017.
- [2] K Cora Ames, Stephen I Ryu, and Krishna V Shenoy. Neural dynamics of reaching following incorrect or absent motor preparation. *Neuron*, 81(2):438–451, 2014.
- [3] Francesco P Battaglia, Gary R Sutherland, and Bruce L McNaughton. Local sensory cues and place cell directionality: additional evidence of prospective coding in the hippocampus. *Journal of Neuroscience*, 24(19):4541–4550, 2004.
- [4] Kevin L Briggman, Henry DI Abarbanel, and William B Kristan. Optical imaging of neuronal populations during decision-making. *Science*, 307(5711):896–901, 2005.
- [5] Bede M Broome, Vivek Jayaraman, and Gilles Laurent. Encoding and decoding of overlapping odor sequences. *Neuron*, 51(4):467–482, 2006.
- [6] Emery N Brown, Loren M Frank, Dengda Tang, Michael C Quirk, and Matthew A Wilson. A statistical paradigm for neural spike train decoding applied to position prediction from ensemble firing patterns of rat hippocampal place cells. *Journal of Neuroscience*, 18(18):7411–7425, 1998.
- [7] Lars Buesing, Jakob H Macke, and Maneesh Sahani. Spectral learning of linear dynamics from generalised-linear observations with application to neural population data. In *Advances in neural information processing systems*, pages 1682–1690, 2012.
- [8] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in beta-vae. *arXiv preprint arXiv:1804.03599*, 2018.
- [9] György Buzsáki. Theta oscillations in the hippocampus. *Neuron*, 33(3):325–340, 2002.
- [10] M Yu Byron, John P Cunningham, Gopal Santhanam, Stephen I Ryu, Krishna V Shenoy, and Maneesh Sahani. Gaussian-process factor analysis for low-dimensional single-trial analysis of neural population activity. In *Advances in neural information processing systems*, pages 1881–1888, 2009.
- [11] Adam J Calhoun, Jonathan W Pillow, and Mala Murthy. Unsupervised identification of the internal states that shape natural behavior. *Nature neuroscience*, 22(12):2040–2049, 2019.
- [12] Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 2610–2620, 2018.

- [13] Mark M Churchland, John P Cunningham, Matthew T Kaufman, Justin D Foster, Paul Nuyujian, Stephen I Ryu, and Krishna V Shenoy. Neural population dynamics during reaching. *Nature*, 487(7405):51–56, 2012.
- [14] Thomas J Davidson, Fabian Kloosterman, and Matthew A Wilson. Hippocampal replay of extended experience. *Neuron*, 63(4):497–507, 2009.
- [15] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- [16] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016.
- [17] Lea Duncker, Gergo Bohner, Julien Boussard, and Maneesh Sahani. Learning interpretable continuous-time models of latent stochastic dynamical systems. *arXiv preprint arXiv:1902.04420*, 2019.
- [18] Lea Duncker and Maneesh Sahani. Temporal alignment and latent gaussian process factor inference in population spike trains. In *Advances in Neural Information Processing Systems*, pages 10445–10455, 2018.
- [19] Alexander S Ecker, Philipp Berens, R James Cotton, Manivannan Subramaniyan, George H Denfield, Cathryn R Cadwell, Stelios M Smirnakis, Matthias Bethge, and Andreas S Tolias. State dependence of noise correlations in macaque primary visual cortex. *Neuron*, 82(1):235–248, 2014.
- [20] David J Foster and Matthew A Wilson. Hippocampal theta sequences. *Hippocampus*, 17(11):1093–1099, 2007.
- [21] Juan A Gallego, Matthew G Perich, Raeed H Chowdhury, Sara A Solla, and Lee E Miller. Long-term stability of cortical population dynamics underlying consistent behavior. *Nature Neuroscience*, pages 1–11, 2020.
- [22] Yuanjun Gao, Evan W Archer, Liam Paninski, and John P Cunningham. Linear dynamical neural population models through nonlinear embeddings. In *Advances in neural information processing systems*, pages 163–171, 2016.
- [23] Apostolos P Georgopoulos, John F Kalaska, Roberto Caminiti, and Joe T Massey. On the relations between the direction of two-dimensional arm movements and cell discharge in primate motor cortex. *Journal of Neuroscience*, 2(11):1527–1537, 1982.
- [24] Apostolos P Georgopoulos, Andrew B Schwartz, and Ronald E Kettner. Neuronal population coding of movement direction. *Science*, 233(4771):1416–1419, 1986.
- [25] Robbe LT Goris, J Anthony Movshon, and Eero P Simoncelli. Partitioning neuronal variability. *Nature neuroscience*, 17(6):858, 2014.
- [26] AD Grosmark, JD Long, and G Buzsáki. Recordings from hippocampal area ca1, pre, during and post novel spatial learning. *crcns. org*, 2016. doi:[10.6080/K0862DC5](https://doi.org/10.6080/K0862DC5).
- [27] Andres D Grosmark and György Buzsáki. Diversity in neural firing dynamics supports both rigid and learned hippocampal sequences. *Science*, 351(6280):1440–1443, 2016.
- [28] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *Iclr*, 2(5):6, 2017.
- [29] Matthew D Hoffman and Matthew J Johnson. Elbo surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference, NIPS*, volume 1, page 2, 2016.
- [30] John Huxter, Neil Burgess, and John O’Keefe. Independent rate and temporal coding in hippocampal pyramidal cells. *Nature*, 425(6960):828–832, 2003.

- [31] Matthew J Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in neural information processing systems*, pages 2946–2954, 2016.
- [32] Jonathan C Kao, Paul Nuyujukian, Stephen I Ryu, Mark M Churchland, John P Cunningham, and Krishna V Shenoy. Single-trial dynamics of motor cortex and their applications to brain-machine interfaces. *Nature communications*, 6(1):1–12, 2015.
- [33] Kenneth Kay, Jason E Chung, Marielena Sosa, Jonathan S Schor, Mattias P Karlsson, Margaret C Larkin, Daniel F Liu, and Loren M Frank. Constant sub-second cycling between representations of possible futures in the hippocampus. *Cell*, 180(3):552–567, 2020.
- [34] Mohammad Reza Keshtkaran and Chethan Pandarinath. Enabling hyperparameter optimization in sequential autoencoders for spiking neural data. In *Advances in Neural Information Processing Systems*, pages 15911–15921, 2019.
- [35] Ilyes Khemakhem, Diederik P Kingma, and Aapo Hyvärinen. Variational autoencoders and nonlinear ica: A unifying framework. *arXiv preprint arXiv:1907.04809*, 2019.
- [36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [37] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [38] Dmitry Kobak, Wieland Brendel, Christos Constantinidis, Claudia E Feierstein, Adam Kepecs, Zachary F Mainen, Xue-Lian Qi, Ranulfo Romo, Naoshige Uchida, and Christian K Machens. Demixed principal component analysis of neural population data. *Elife*, 5:e10989, 2016.
- [39] Vernon Lawhern, Wei Wu, Nicholas Hatsopoulos, and Liam Paninski. Population decoding of motor cortical activity using a generalized linear model with hidden states. *Journal of neuroscience methods*, 189(2):267–280, 2010.
- [40] Mikhail A Lebedev, Alexei Ossadtchi, Nil Adell Mill, Núria Armengol Urpí, Maria R Cervera, and Miguel AL Nicolelis. Analysis of neuronal ensemble activity reveals the pitfalls and shortcomings of rotation dynamics. *Scientific Reports*, 9(1):1–14, 2019.
- [41] Shen Li, Bryan Hooi, and Gim Hee Lee. Identifying through flows for recovering latent representations. *arXiv preprint arXiv:1909.12555*, 2019.
- [42] I-Chun Lin, Michael Okun, Matteo Carandini, and Kenneth D Harris. The nature of shared cortical variability. *Neuron*, 87(3):644–656, 2015.
- [43] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. *arXiv preprint arXiv:1811.12359*, 2018.
- [44] Jakob H Macke, Lars Buesing, John P Cunningham, M Yu Byron, Krishna V Shenoy, and Maneesh Sahani. Empirical models of spiking in neural populations. In *Advances in neural information processing systems*, pages 1350–1358, 2011.
- [45] Valerio Mante, David Sussillo, Krishna V Shenoy, and William T Newsome. Context-dependent computation by recurrent dynamics in prefrontal cortex. *nature*, 503(7474):78–84, 2013.
- [46] Ofer Mazor and Gilles Laurent. Transient dynamics versus fixed points in odor representations by locust antennal lobe projection neurons. *Neuron*, 48(4):661–673, 2005.
- [47] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- [48] Zaneta Navratilova, Lan T Hoang, C Daniela Schwindel, Masami Tatsuno, and Bruce L McNaughton. Experience-dependent firing rate remapping generates directional selectivity in hippocampal place cells. *Frontiers in neural circuits*, 6:6, 2012.

- [49] John O’Keefe. Place units in the hippocampus of the freely moving rat. *Experimental neurology*, 51(1):78–109, 1976.
- [50] John O’Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence from unit activity in the freely-moving rat. *Brain research*, 1971.
- [51] Mike W Oram, Peter Földiák, David I Perrett, and Frank Sengpiel. Theideal homunculus’: decoding neural population signals. *Trends in neurosciences*, 21(6):259–265, 1998.
- [52] Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky, Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg, et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature methods*, 15(10):805–815, 2018.
- [53] Liam Paninski, Shy Shoham, Matthew R Fellows, Nicholas G Hatsopoulos, and John P Donoghue. Superlinear population encoding of dynamic hand trajectory in primary motor cortex. *Journal of Neuroscience*, 24(39):8551–8561, 2004.
- [54] Il Memming Park, Miriam LR Meister, Alexander C Huk, and Jonathan W Pillow. Encoding and decoding in parietal cortex during sensorimotor decision-making. *Nature neuroscience*, 17(10):1395, 2014.
- [55] Eva Pastalkova, Vladimir Itskov, Asohan Amarasingham, and György Buzsáki. Internally generated cell assembly sequences in the rat hippocampus. *Science*, 321(5894):1322–1327, 2008.
- [56] Jonathan W Pillow, Jonathon Shlens, Liam Paninski, Alexander Sher, Alan M Litke, EJ Chichilnisky, and Eero P Simoncelli. Spatio-temporal correlations and visual signalling in a complete neuronal population. *Nature*, 454(7207):995–999, 2008.
- [57] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- [58] Patrick T Sadtler, Kristin M Quick, Matthew D Golub, Steven M Chase, Stephen I Ryu, Elizabeth C Tyler-Kabara, M Yu Byron, and Aaron P Batista. Neural constraints on learning. *Nature*, 512(7515):423–426, 2014.
- [59] Terence David Sanger. Probability density estimation for the interpretation of neural population codes. *Journal of neurophysiology*, 76(4):2790–2793, 1996.
- [60] Andrew B Schwartz, Ronald E Kettner, and Apostolos P Georgopoulos. Primate motor cortex and free arm movements to visual targets in three-dimensional space. i. relations between single cell discharge and direction of movement. *Journal of Neuroscience*, 8(8):2913–2927, 1988.
- [61] Krishna V Shenoy, Maneesh Sahani, and Mark M Churchland. Cortical control of arm movements: a dynamical systems perspective. *Annual review of neuroscience*, 36:337–359, 2013.
- [62] William E Skaggs, Bruce L McNaughton, Matthew A Wilson, and Carol A Barnes. Theta phase precession in hippocampal neuronal populations and the compression of temporal sequences. *Hippocampus*, 6(2):149–172, 1996.
- [63] Herman P Snippe. Parameter extraction from population codes: A critical assessment. *Neural Computation*, 8(3):511–529, 1996.
- [64] Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ica with general incompressible-flow networks (gin). *arXiv preprint arXiv:2001.04872*, 2020.
- [65] Lakshminarayanan Srinivasan, Uri T Eden, Alan S Willsky, and Emery N Brown. A state-space analysis for reconstruction of goal-directed movements using neural signals. *Neural computation*, 18(10):2465–2494, 2006.
- [66] Mark Stopfer, Vivek Jayaraman, and Gilles Laurent. Intensity versus identity coding in an olfactory system. *Neuron*, 39(6):991–1004, 2003.

- [67] Wilson Truccolo, Uri T Eden, Matthew R Fellows, John P Donoghue, and Emery N Brown. A point process framework for relating neural spiking activity to spiking history, neural ensemble, and extrinsic covariate effects. *Journal of neurophysiology*, 93(2):1074–1089, 2005.
- [68] Xue-Xin Wei, Ding Zhou, Andres Grosmark, Zaki Ajabi, Fraser Sparks, Pengcheng Zhou, Mark Brandon, Attila Losonczy, and Liam Paninski. A zero-inflated gamma model for post-deconvolved calcium imaging traces. *bioRxiv*, page 637652, 2019.
- [69] IQ Whishaw and C Hippocampal Vanderwolf. Hippocampal eeg and behavior: change in amplitude and frequency of rsa (theta rhythm) associated with spontaneous and learned movement patterns in rats and cats. *Behavioral biology*, 8(4):461–484, 1973.
- [70] Matthew R Whiteway and Daniel A Butts. The quest for interpretable models of neural population activity. *Current opinion in neurobiology*, 58:86–93, 2019.
- [71] Anqi Wu, Nicholas A Roy, Stephen Keeley, and Jonathan W Pillow. Gaussian process based nonlinear latent structure discovery in multivariate spike train data. In *Advances in neural information processing systems*, pages 3496–3505, 2017.
- [72] Wei Wu, Jayant E Kulkarni, Nicholas G Hatsopoulos, and Liam Paninski. Neural decoding of hand motion using a linear state-space model with hidden states. *IEEE Transactions on neural systems and rehabilitation engineering*, 17(4):370–378, 2009.
- [73] Kechen Zhang, Iris Ginzburg, Bruce L McNaughton, and Terrence J Sejnowski. Interpreting neuronal population activity by reconstruction: unified framework with application to hippocampal place cells. *Journal of neurophysiology*, 79(2):1017–1044, 1998.
- [74] Yuan Zhao and Il Memming Park. Variational latent gaussian process for recovering single-trial dynamics from population spike trains. *Neural computation*, 29(5):1293–1316, 2017.

Supplementary Information (SI)

A Proof of identifiability for pi-VAE

Theorem 1. Assume that we observe data sampled from pi-VAE model defined according to equation 1,2 with Poisson noise and parameters $\theta = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$. Assume the following holds:

- i) The firing rate function \mathbf{f} in equation 1 is injective.
- ii) The sufficient statistics $T_{i,j}$ in 2 are differentiable almost everywhere, and their derivatives $T'_{i,j}$ are nonzero almost everywhere for $1 \leq i \leq m, 1 \leq j \leq k$.
- iii) There exists $mk + 1$ distinct points $\mathbf{u}^0, \dots, \mathbf{u}^{mk+1}$ such that the matrix

$$L = (\boldsymbol{\lambda}(\mathbf{u}^1) - \boldsymbol{\lambda}(\mathbf{u}^0), \dots, \boldsymbol{\lambda}(\mathbf{u}^{mk}) - \boldsymbol{\lambda}(\mathbf{u}^0))$$

of size $mk \times mk$ is invertible, then the pi-VAE model is identifiable up to \sim .

Proof. [35] has proved that the Bernoulli observation model is identifiable under the same set of assumptions. For Poisson observations with mean firing rate as λ , we can transform it to Bernoulli observations with parameter $p = 1 - \exp(-\lambda)$ by keeping the zeros and treating the positive values as ones. Because the Bernoulli model is identifiable, the Poisson model is also identifiable. \square

B Network architecture

For both the generative models in pi-VAE and VAE, we used the following strategy to parameterize $\mathbf{f}(\cdot)$ which maps the m -dimensional latent \mathbf{z} to the mean firing rate of n -dimensional Poisson observations. We first mapped the $\mathbf{z}_{1:m}$ to the concatenation of $\mathbf{z}_{1:m}$ and $t(\mathbf{z}_{1:m})$, where $t(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}^{n-m}$ is parameterized by a feed-forward neural network with a linear output and 2 hidden layers, each containing $\lfloor n/4 \rfloor$ nodes with ReLU activation function. Then we applied two GIN blocks. Same as [64], we defined the affine coupling function as the concatenation of the scale function s and the translation function t , computed together for efficiency, applied two affine coupling functions per GIN block, and randomly permuted the input before passing it through each GIN block. We defined both s, t in GIN block as mapping: $\mathbb{R}^{\lfloor n/2 \rfloor} \rightarrow \mathbb{R}^{n-\lfloor n/2 \rfloor}$. The scale function s is passed through a clamping function $0.1 \tanh(s)$, which limits the output to the range $(-0.1, 0.1)$. For affine coupling function, we have a linear output layer with 2 hidden layers, each containing $\lfloor n/4 \rfloor$ nodes with ReLU activation function.

We modeled the prior $p_{T,\lambda}(\mathbf{z}|\mathbf{u})$ in pi-VAE as independent Gaussian distribution. The natural parameters $\lambda_{i,j}$ are the Gaussian means and variances. For discrete \mathbf{u} , we used different values of the mean and variance for different labels. For continuous \mathbf{u} , we parameterized the mean and spectrum decomposition of variance together by a feed-forward neural network with a linear output layer and 2 hidden layers, each containing 20 nodes with tanh activation function (the mean and variance share the 2 hidden layers). For the cases of mixed discrete and continuous labels \mathbf{u} , we encoded the discrete labels with a one-hot vector, and mapped it together with the continuous components to the mean and spectrum decomposition of variance using feed-forward neural network as described in the continuous case.

For the recognition model in pi-VAE and VAE, we used $q_\phi(\mathbf{z}|\mathbf{x})$ as independent Gaussian distribution, and parameterized the mean and the spectrum decomposition of the variance separately using feed-forward neural network with a linear output layer and 2 hidden layers, each containing 60 nodes with tanh activation function.

Code implementing the algorithms is available at <https://github.com/zhd96/pi-vae>.

C Synthetic data simulations

To generate firing rate of the Poisson process from simulated latent \mathbf{z} , we first padded \mathbf{z} with $n-m$ zeros, then applied 4 RealNVP blocks, each containing 2 affine coupling functions with the same structure as defined in section B except that s does not need to have sum equal to 0 here, and we used $\lfloor n/2 \rfloor$ nodes for each hidden layer.

For discrete label simulation shown in Fig. 2a-d, we simulated 10^4 observations, and split them into training, validation, test data (80%, 10%, 10% respectively). We set the batch size to be 200 during training, and trained for 600 epochs. For the continuous label simulation shown in Fig. 2e-h, we simulated 1.5×10^4 observations. The training-validation-test split is the same as discrete label simulation. We set batch size as 300, and trained for 1000 epochs.

D Monkey reaching data: session 2

For each reaching direction, there are ~ 35 trials (see Fig. S1a). We analyzed 211 neurons from PMd area, and focused on the reaching period from go cue (defined as $t = 0$) to the end, which typically last for ~ 1 second. We binned the ensemble spike activities into 50ms bins. We randomly split the dataset into 34 batches, where each batch contains at least one trial from each direction. We randomly took 28 batches as training data, 3 batches as validation data and 3 batches as test data. Similar to Session 1, We fit 4-dimensional latent models to the data based on pi-VAE and VAE respectively. Results are shown in Fig. S1.

E Alternative methods

E.1 Monkey reaching data

For supervised UMAP², we set the reaching directions as labels, and embedded the high dimensional spike count data into a 4-dimensional latent space. Other parameters were set to be the default values. For PfLDS, we implemented the algorithm on our own using the same neural network architecture as in [22] (the original code provided by the authors of that paper depends on Python Theano library, which has not been maintained for a while). We assumed a 4-dimensional latent space and Poisson observation model. We set the learning rate as 2.5×10^{-4} and trained for 1500 epochs. Each batch consisted of a single trial. The training, validation and test sets had 184, 16, 17 trials respectively. For LFADS³, we assumed a 4-dimensional latent space along with a Poisson observation model. We applied two versions of the model to the data, with the reaching direction as an additional input and without this input. Other parameters were set to be the default values. We pre-processed the data by discarding all the trials less than 1 second and trimming longer trials to make them 1 second long. Each batch consisted of a single trial. The training, validation and test sets had 177, 16, 16 trials respectively. For demixed PCA⁴, we pre-processed the data to make each reaching direction had the same number of trials and each trial had the same length (*i.e.*, 1 second). We took time and stimulus as labels. We used 2 sets of components, each containing time, stimulus, as well as time and stimulus mixing components. Other parameters (eg. regularizer) were set to be the default values.

E.2 Hippocampus data

For supervised UMAP, we used 2-dimensional latent variables model with rat's locations as labels. Other parameters were set as default. For PCA, we used two principal components. For "PCA after LDA", we first applied LDA with rat's running directions as response and neural activities as predictors, and identified the 1-dimensional linear boundary which could separate the neural activities of two directions most. Then we projected the neural activities on this boundary using linear regression, then applied PCA with 2 principal components on the residuals. We found that the resulting latents were all less interpretable than pi-VAE (Fig. S4), with no dimension directly representing the rat's location. Also the rhythmic-like fluctuations spanned across dimensions, rather than concentrated in one dimension (not shown).

²<https://umap-learn.readthedocs.io/en/latest/>

³<https://github.com/lfads/models>

⁴<https://github.com/machenslab/dPCA>

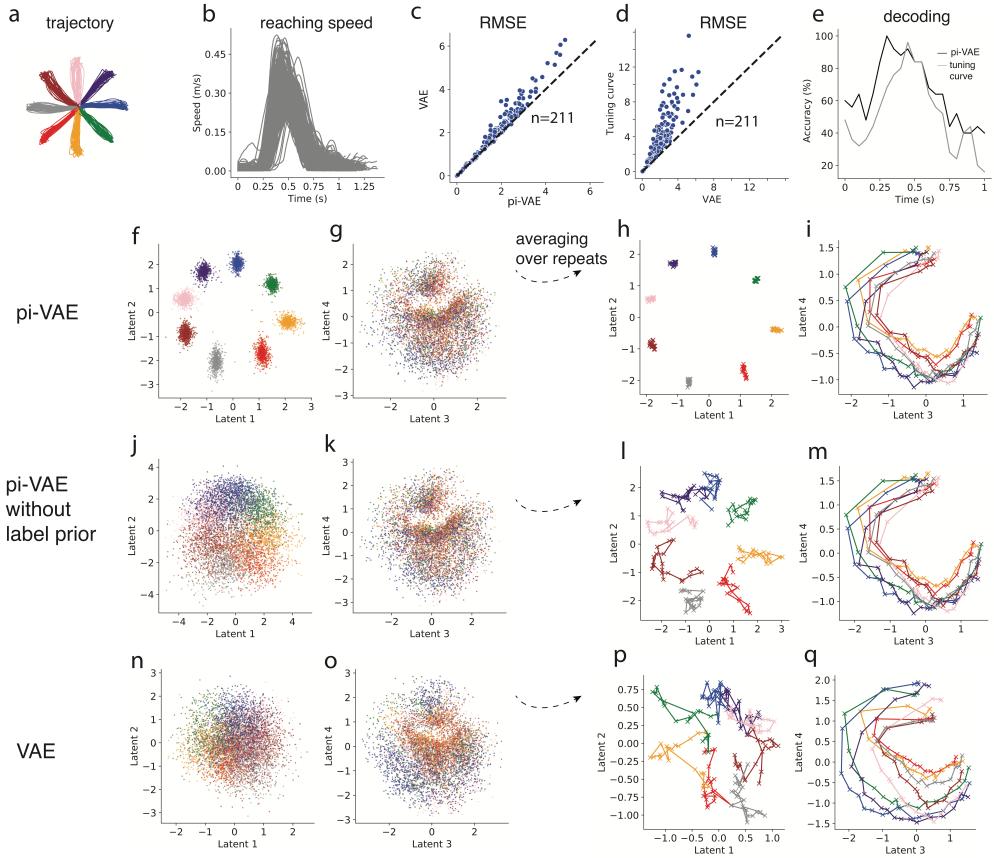


Figure S1: Results on monkey reaching data (Session 2). These results are similar to those obtained from Session 1 as reported in the main text. (a) The macaque’s reaching trajectories for 8 directions labeled by different colors. (b) The reaching speed of the macaque for each trial. (c,d) Scatter plots of RMSE of fitted rate ($n = 211$ neurons) for comparing pi-VAE and VAE, as well as VAE and tuning curve. (e) Decoding accuracy as function of time on test data by pi-VAE and tuning curve model. (f,g) Inferred latent based on pi-VAE, i.e.,mean of $q(\mathbf{z}|\mathbf{x}, \mathbf{u})$. (h,i) Inferred latent from pi-VAE averaged over repeats from the same reaching direction. (j,k) Mean of $q(\mathbf{z}|\mathbf{x})$ from pi-VAE. (l,m) Mean of $q(\mathbf{z}|\mathbf{x})$ by pi-VAE averaging over repeats from the same reaching direction. (n-q) Similar to (f-i) for VAE.

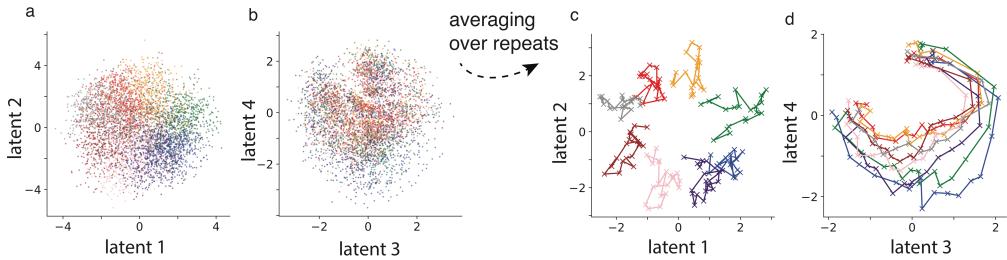


Figure S2: Related to Fig. 3, on reaching data. Inferred latent without label prior using pi-VAE still are still highly structured and interpretable. The first two dimensions carry information about the reaching direction, while the third and fourth dimension mainly captures the dynamics over the time course of a trial. (a,b) Mean of $q(\mathbf{z}|\mathbf{x})$ from pi-VAE. (c,d) Mean of $q(\mathbf{z}|\mathbf{x})$ by pi-VAE averaging over repeats from the same reaching direction.

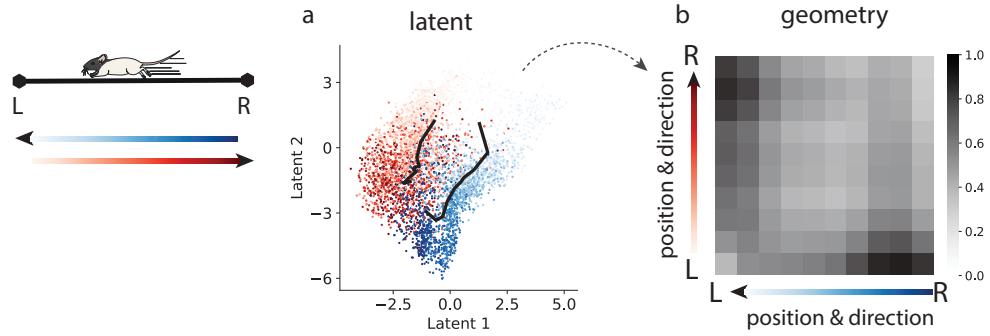


Figure S3: Related to Fig. 4. Results on Hippocampus CA1 data. Inferred latents without the label prior using pi-VAE still exhibit clear structure, with the latent geometry respecting the geometry of the track. (a) Mean of $q(\mathbf{z}|\mathbf{x})$ from pi-VAE. Two directions are color-coded by red and blue, and positions are coded by color saturation. Black lines represent the mean of the latent states corresponding to position on the track for two directions. (b) The distance between pairs of points from the two black lines is computed to quantify the latent geometry.

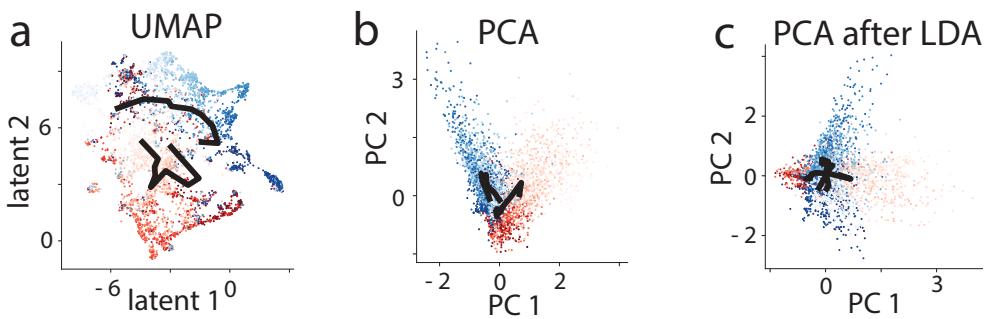


Figure S4: Hippocampus data: results from several alternative methods. a) UMAP. b) PCA. c) PCA after Linear Discriminant analysis (LDA). Notice that these methods recovered more entangled representation compared to pi-VAE.