

BRAID: INPUT-DRIVEN NONLINEAR DYNAMICAL MODELING OF NEURAL-BEHAVIORAL DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural populations exhibit complex recurrent structures that drive behavior, while continuously receiving and integrating external inputs from sensory stimuli, upstream regions, and neurostimulation. However, neural populations are often modeled as autonomous dynamical systems, with little consideration given to the influence of external inputs that shape the population activity and behavioral outcomes. Here, we introduce BRAID, a deep learning framework that models nonlinear neural dynamics underlying behavior while explicitly incorporating any measured external inputs. Our method disentangles intrinsic recurrent neural population dynamics from the effects of inputs by including a forecasting objective within input-driven recurrent neural networks. BRAID further prioritizes the learning of intrinsic dynamics that are related to a behavior of interest by using a multi-stage optimization scheme. We validate BRAID with nonlinear simulations, showing that it can accurately learn the intrinsic dynamics shared between neural and behavioral modalities. We then apply BRAID to motor cortical activity recorded during a motor task and demonstrate that our method more accurately fits the neural-behavioral data by incorporating measured sensory stimuli into the model and improves the forecasting of neural-behavioral data compared with various baseline methods, whether input-driven or not.

1 INTRODUCTION

Understanding the relationship between neural activity and behavior is a critical goal in neuroscience and neurotechnology. Neural activity and its temporal structure, or “*dynamics*” during a behavior, are formed by the interplay between (1) the recurrent networks within a brain area, i.e., *intrinsic dynamics*, and the (2) temporally-structured inputs it receives during the behavior (Remington et al., 2018; Vyas et al., 2020). A brain population may receive inputs from measurable sources such as sensory stimuli, electrical/optogenetic neurostimulation (Buonomano & Maass, 2009; Seely et al., 2016; Susilaradeya et al., 2019; Sauerbrei et al., 2020; Shenoy & Kao, 2021; Vahidi et al., 2024), as well as from other upstream brain areas (Sauerbrei et al., 2020; Shenoy & Kao, 2021), which could be included in multi-regional recordings (Jun et al., 2017; Steinmetz et al., 2019). However, even easily measurable external inputs (e.g., sensory stimuli) are often not explicitly considered when modeling neural-behavioral activity, which can lead to a conflation of intrinsic and input-driven contributions, creating challenges for interpretation (Seely et al., 2016; Sauerbrei et al., 2020; Vahidi et al., 2024). Beyond disentangling intrinsic dynamics from input dynamics, incorporating measured inputs into models can also enhance the behavior decoding performance in neurotechnologies such as stimulation-based closed-loop controllers (Yang et al., 2021).

Another important challenge is to disentangle neural dynamics that are relevant to a specific behavior from other neural dynamics, and to prioritize the former. This is critical because the majority of neural variance may not be relevant to the behavior of interest (Churchland et al., 2012; Mante et al., 2013; Kobak et al., 2016; Allen et al., 2019; Engel & Steinmetz, 2019; Stringer et al., 2019; Sani et al., 2021). While most prior works have used unsupervised approaches when modeling neural activity as latent variable dynamical systems (Aghagolzadeh & Truccolo, 2015; Gao et al., 2016; Wu et al., 2017; Pandarinath et al., 2018; Hernandez et al., 2020; Rutten et al., 2020; Kim et al., 2021), more recent works have shown improved learning of behaviorally relevant neural dynamics by using behavior data during learning in a supervised manner (Hurwitz et al., 2021; Sani et al., 2021; Kramer et al., 2022; Gondur et al., 2024; Vahidi et al., 2024; Sani et al., 2024).

054 Yet another challenge is posed by the nonlinearities in neural-behavioral data. While linear models
055 have been extremely effective in approximating neural dynamics (Hastie et al., 2009; Churchland
056 et al., 2012; Mante et al., 2013; Cunningham & Yu, 2014; Kao et al., 2015; Kobak et al., 2016;
057 Abbaspourazad et al., 2021; Sani et al., 2021), they may require higher dimensional latent repre-
058 sentations compared to nonlinear models (Yang et al., 2019; Nozari et al., 2024; Sani et al., 2024),
059 and do not provide interpretability for nonlinear dynamical phenomena such as multi-stable fixed
060 points and limit cycles (Kim et al., 2021; Durstewitz et al., 2023). Moreover, unlike linear models,
061 for nonlinear models the relationship between the intrinsic dynamics and an inference model that is
062 fitted to estimate the latent states from observations is not analytically known, posing a challenge
063 for studying intrinsic dynamics (see section 3.1).

064 Here, we address all aforementioned challenges by introducing Behaviorally Relevant Analysis of
065 Intrinsic Dynamics (BRAID), a new method with the following key contributions. *First*, BRAID
066 captures complex nonlinear structures in neural-behavioral-input data, offering greater expressiv-
067 ity than linear methods. *Second*, by optimizing multi-step-ahead forecasts of neural-behavior data,
068 BRAID simultaneously learns two representations for neural dynamics: the predictor and the gener-
069 ative form representations (see section 3.1), the latter of which describes intrinsic dynamics. *Third*,
070 by explicitly modeling the influence of measured inputs, BRAID disentangles their dynamics from
071 intrinsic dynamics, to more closely reflect the neuronal networks within the recorded brain region.
072 *Fourth*, we introduce a multi-stage learning framework that dissociates and prioritizes the learning
073 of intrinsic behaviorally relevant neural dynamics, while considering measured inputs (see section
074 3.2). *Fifth*, we introduce additional preprocessing and post-hoc learning stages that allow behavior-
075 specific dynamics to be dissociated from behaviorally relevant neural dynamics (see section 3.3).

076 We validate BRAID in multiple simulated datasets with distinct nonlinear structures and show its
077 capability to accurately learn the underlying nonlinear model, resulting in an interpretable represen-
078 tation of intrinsic dynamics. We then apply our method to electrophysiological data recorded from a
079 non-human-primate (NHP) performing sequential reaches (O’Doherty et al., 2017). Our results indi-
080 cate that accounting for both nonlinearity and sensory inputs improves neural-behavioral prediction
081 suggesting a more accurate representation of intrinsic behaviorally relevant neural dynamics.

082 2 RELATED WORK

083 Our work addresses multiple problems simultaneously, which makes it related to various methods
084 that tackle a subset of these problems. A summary of related methods is provided in table 1.

085 First, a key ability of BRAID is to incorporate measured inputs to disentangle intrinsic dynamics
086 from input dynamics. Other nonlinear modeling methods, including those based on deep learning
087 (Gao et al., 2016; Sussillo et al., 2016; Wu et al., 2017; Pandarinath et al., 2018; Rutten et al., 2020;
088 Hurwitz et al., 2021; Kim et al., 2021; Abbaspourazad et al., 2024; Sani et al., 2024), have not ad-
089 dressed the problem of modeling measured external inputs and their impact on neural-behavioral
090 data. As demonstrated by Vahidi et al. (2024), not considering external inputs can lead to the dy-
091 namics of these inputs being misinterpreted as intrinsic neural dynamics. To overcome this chal-
092 lenge, Vahidi et al. (2024) introduce a linear dynamical modeling method, termed IPSID, which
093 explicitly incorporates measured external inputs into the model. However, IPSID is an analytical,
094 projection-based and strictly linear method that cannot capture any nonlinearities. By incorporating
095 the strengths of this linear modeling work into BRAID, we can account for measured external inputs
096 and dissociate their dynamics while allowing every element of the model to be nonlinear. We use
097 IPSID as a key baseline to show the benefit of enabling nonlinearity in our method (see appendix
098 A.2 for details). We also show the results for the special case of setting all model element as linear
099 in our method (referred to as linear BRAID), which fits in a linear model similar to that of IPSID.

100 Second, a key capability of BRAID is that it dissociates behaviorally relevant neural dynamics into
101 a distinct part of the latent states and prioritizes their learning, while also being able to learn neural-
102 specific and behavior-specific dynamics using additional latent states. Among prior works, two
103 recent nonlinear methods termed DPAD (Sani et al., 2024) and TNDM (Hurwitz et al., 2021) aim to
104 dissociate behaviorally relevant dynamics from other neural dynamics, but neither method dissoci-
105 ates the third category of dynamics, i.e., the behavior-specific dynamics. More importantly, neither
106 DPAD nor TNDM incorporates external inputs into the model to dissociate intrinsic dynamics from
107 input dynamics. Finally, DPAD learns models based on 1-step-ahead prediction of neural-behavioral

108 data and does not explicitly learn the intrinsic dynamics, whereas BRAID adds m -step-ahead pre-
 109 dictions into the loss to optimize forecasting and also explicitly learns a generative representation
 110 of intrinsic dynamics. TNDM on the other hand is a sequential autoencoders (similar to LFADS,
 111 Pandarinath et al., 2018), i.e., it optimizes reconstruction of a window of data after ingesting the
 112 entire window as input. We compared our results with both DPAD and TNDM, although DPAD’s
 113 architecture is closer to ours. In fact, the comparisons with DPAD can also be thought of as ablation
 114 studies that show the benefit of incorporating external inputs and forecasting in our method.

115 Third, we learn behaviorally relevant neural dynamics, or in other words the shared neural-
 116 behavioral dynamics, in an initial optimization focused on learning these dynamics, while leaving
 117 the learning of other neural dynamics to a separate subsequent optimization. This approach, priori-
 118 tizes behaviorally relevant neural dynamics in the sense that we can fit models with low dimensional
 119 latent states that are purely focused on these dynamics. Besides DPAD, a few other works, includ-
 120 ing TNDM, propose nonlinear approaches for learning dynamics shared between two modalities
 121 (Hurwitz et al., 2021; Kramer et al., 2022; Gondur et al., 2024). However, these works, use a com-
 122 bined loss to optimize the reconstruction of both modalities in the same optimization. While this
 123 approach can capture the dynamics shared between modalities, it does not prioritize them over dy-
 124 namics specific to either modality (Sani et al., 2024). Moreover, most multi-modal approaches do
 125 not model the effect of external inputs (Hurwitz et al., 2021; Gondur et al., 2024). One multi-modal
 126 framework, termed mmPLRNN (Kramer et al., 2022), which models dynamics of two modalities
 127 with a piecewise-linear RNN (see appendix A.2 for details), supports modeling the effect of exter-
 128 nal inputs, although this capability was not demonstrated in Kramer et al. (2022). Nevertheless,
 129 we include comparisons with mmPLRNN with input as one of our baselines. Finally, as another
 130 ablation study to assess the importance of prioritization of behaviorally relevant dynamics, we also
 131 implement an unsupervised version of BRAID, termed U-BRAID, that removes the behaviorally
 132 relevant optimization step and instead learns all neural dynamics in one optimization step while still
 incorporating external inputs into the model (see section 3 and appendix A.2 for details).

133 Most other prior nonlinear methods only consider neural signals during modeling without consider-
 134 ing behavior or external inputs (Gao et al., 2016; Pandarinath et al., 2018; Hernandez et al., 2020;
 135 Rutten et al., 2020; Kim et al., 2021) or do not use dynamic models (Zhou & Wei, 2020; Schneider
 136 et al., 2023) and thus are vastly different from our method. Nevertheless, we include comparisons
 137 with LFADS (Pandarinath et al., 2018) and CEBRA (Schneider et al., 2023) as additional baselines.
 138 We list the differences of some of these methods with our method in table 1.

140 Table 1: Related works (see Discussion). ELBO: evidence lower bound, LL: log-likelihood.

Method	Nonlinear	Prioritize behaviorally relevant	Dissociate non-neural	Dissociate intrinsic	Training objective
IPSID	✗	✓	✓	✓	Projection-based
TNDM	✓	✓	✗	✗	Multi-modal ELBO
LFADS	✓	✗	✗	✗	ELBO
mmPLRNN	✓	✗	✗	✓	Multi-modal ELBO
DPAD	✓	✓	✗	✗	1-step-ahead LL
CEBRA	✓	✓	✗	✗	Contrastive
BRAID	✓	✓	✓	✓	m -step-ahead LL

152 3 METHODS

153 3.1 BRAID MODEL

154 We model the neural activity ($\mathbf{y}_k \in \mathbb{R}^{n_y}$) and behavior ($\mathbf{z}_k \in \mathbb{R}^{n_z}$) as observations of a nonlinear
 155 dynamical system (with latent states $\mathbf{x}_k^s \in \mathbb{R}^{n_x}$) that have some intrinsic dynamics and are driven
 156 both by measured inputs ($\mathbf{u}_k \in \mathbb{R}^{n_u}$) as well as unmeasured inputs and/or noises ($\mathbf{w}_k \in \mathbb{R}^{n_x}$). The
 157 following model describes the dynamical system

$$\begin{cases} \mathbf{x}_{k+1}^s &= \mathbf{A}_{fw}(\mathbf{x}_k^s) + \mathbf{K}_{fw}(\mathbf{u}_k) + \mathbf{w}_k \\ \mathbf{y}_k &= \mathbf{C}_y(\mathbf{x}_k^s, \mathbf{u}_k) + \mathbf{v}_k \\ \mathbf{z}_k &= \mathbf{C}_z(\mathbf{x}_k^s, \mathbf{u}_k) + \boldsymbol{\epsilon}_k \end{cases} \quad (1)$$

where $\mathbf{v}_k \in \mathbb{R}^{n_y}$ and $\epsilon_k \in \mathbb{R}^{n_z}$ are observation noises. Given this dynamical system, one can recursively infer the latent state from neural observations \mathbf{y}_k using an RNN as follows

$$\mathbf{x}_{k+1|k} = \mathbf{A}(\mathbf{x}_{k|k-1}) + \mathbf{K}(\mathbf{y}_k, \mathbf{u}_k) \quad (2)$$

where $\mathbf{x}_{k+1|k}$ (or simply \mathbf{x}_{k+1}) is defined as the inferred latent state based on $\{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ and $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$. Given the latent nature of the states, even when inference is optimal (e.g., in a Kalman filter), the inferred states will *not* be equal to the internal states \mathbf{x}_k^s in equation 1 (Katayama, 2006), which is why we use different notations for the states in equations 1 and 2. More importantly, note that \mathbf{A} and \mathbf{K} in equation 2 are distinct from \mathbf{A}_{fw} and \mathbf{K}_{fw} in equation 1. This is because \mathbf{A} and \mathbf{K} represent the “predictor form” representation of dynamics, describing how the inferred latent state recursively evolves over time as samples of \mathbf{y}_k and \mathbf{u}_k are observed, whereas \mathbf{A}_{fw} and \mathbf{K}_{fw} represent the “generative form” representation of dynamics that describe how the latent states themselves evolve, purely based on their *intrinsic dynamics* – so \mathbf{A}_{fw} is what ultimately describes the intrinsic dynamics. For linear systems, there is an analytical bidirectional relationship between predictor and generative form representations (defined by the Kalman filter, see Katayama, 2006), whereas for nonlinear systems in general, this relationship is not known. Thus, we devise an approach that allows us to learn both representations of dynamics from data.

Critically, to predict the latent state (or neural or behavioral data) multiple (for $m > 1$) steps into the future using only new observations from the input \mathbf{u}_k , we would need to propagate the latent state ahead according to its *intrinsic* dynamics, i.e., the “generative form” representation of dynamics, as

$$\mathbf{x}_{k+m|k} = \mathbf{A}_{fw}(\mathbf{x}_{k+m-1|k}) + \mathbf{K}_{fw}(\mathbf{u}_{k+m-1}) \quad (3)$$

where $\mathbf{x}_{k+m|k}$ denotes the latent state at time step $k + m$, generated given $\{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ and $\{\mathbf{u}_1, \dots, \mathbf{u}_{k+m-1}\}$. Note that for $m = 2$, the right hand side of equation 3 would have $\mathbf{x}_{k+1|k}$, which is given by equation 2. Thus, m -step-ahead inference of the latent state engages both predictor and generative form representation of the dynamics via equations 2 and 3, respectively. As an alternative interpretation, the m -step-ahead prediction of the latent state (or neural or behavioral data), for $m > 1$, involves two RNNs operating in complementary fashion (figure 1b):

1. The first RNN (*RNN*, parameterized by \mathbf{A} and \mathbf{K}) takes in neural and input time series and recursively estimates the 1-step-ahead prediction ($\mathbf{x}_{k|k-1}$).
2. The second RNN (*RNN_{fw}*, parameterized by \mathbf{A}_{fw} and \mathbf{K}_{fw}) takes in the 1-step-ahead predicted state from the first RNN and propagates it $m - 1$ additional steps ahead according to the intrinsic latent dynamics of the model, to get the m -step-ahead predictions ($\mathbf{x}_{k+m-1|k-1}$, for $m > 1$).

Overall, the BRAID model is comprised of six distinct transformations: $\mathbf{A}(\cdot)$, $\mathbf{A}_{fw}(\cdot)$, $\mathbf{K}(\cdot)$, $\mathbf{K}_{fw}(\cdot)$, $\mathbf{C}_z(\cdot)$, and $\mathbf{C}_y(\cdot)$. $\mathbf{A}/\mathbf{A}_{fw}$ describe predictor/generative form recursions of the latent state. $\mathbf{K}/\mathbf{K}_{fw}$ describe predictor/generative form encoders. \mathbf{C}_z and \mathbf{C}_y describe behavior and neural decoders. We implement these six transformations as multi-layer perceptrons (MLPs) with arbitrary user-specified number of units and hidden layers. As a special case, any (or all) of these mappings can be replaced by a linear mapping (i.e., an MLP with no hidden layer and a linear activation).

We learn the parameters specifying all six transformations of the model by optimizing a weighted sum of m -step-ahead neural-behavioral prediction errors (for $m \in [m_1, m_2, \dots, m_L]$) as our losses

$$\begin{aligned} L_z &= \sum_{i=1}^L \alpha_{z_{m_i}} \text{MSE}(\mathbf{z}_{k+m_i}, \mathbf{C}_z(\mathbf{x}_{k+m_i|k}, \mathbf{u}_{k+m_i})) \\ L_y &= \sum_{i=1}^L \alpha_{y_{m_i}} \text{MSE}(\mathbf{y}_{k+m_i}, \mathbf{C}_y(\mathbf{x}_{k+m_i|k}, \mathbf{u}_{k+m_i})) \end{aligned} \quad (4)$$

where $\text{MSE}(\cdot)$ indicates the mean-squared error loss, L denotes the number of steps ahead simultaneously included in the loss, and $\alpha_{z_{m_i}}$ and $\alpha_{y_{m_i}}$ denote the weights used in the sum. In this work, we always set $\alpha_{z_{m_i}}$ and $\alpha_{y_{m_i}}$ to 1. Moreover, although the decoders in BRAID can optionally take both the latent state and the external input \mathbf{u}_k (lines 2-3 of equation 1), in our real data analyses we do not provide \mathbf{u}_k to decoders and generate predictions only based on the latent states.

3.2 PRIORITIZATION OF BEHAVIORALLY RELEVANT OVER OTHER NEURAL DYNAMICS

To dissociate behaviorally relevant neural dynamics from other neural dynamics and prioritize the former, we break the latent state \mathbf{x}_k into two sections ($\mathbf{x}_k^{(1)}$ and $\mathbf{x}_k^{(2)}$) and learn these two sections in two learning stages. We denote the model parameters associated with each model section using

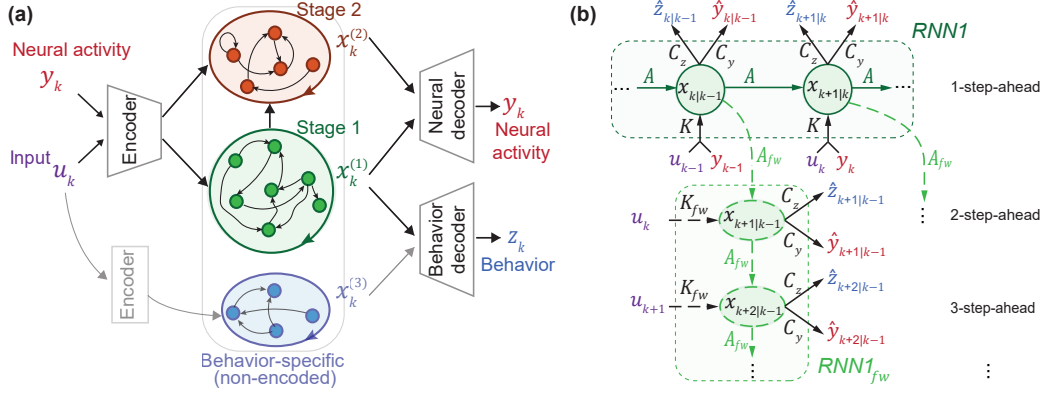


Figure 1: **BRAID model architecture.** (a) BRAID dissociates the dynamics of neural-behavioral data into three latent states $x_k^{(1)}$, $x_k^{(2)}$, and $x_k^{(3)}$: 1) the dynamics shared between neural and behavioral modalities (learned in stage 1 by $RNN1$ and $RNN1_{fw}$), 2) any remaining dynamics private to neural activity (learned in stage 2 by $RNN2$ and $RNN2_{fw}$), and 3) input-driven, behavior-specific dynamics not encoded in neural activity (learned per section 3.3 by $RNN3$ and $RNN3_{fw}$). (b) For each latent state, we simultaneously learn a predictor and generative form representation of the dynamics (denoted by A and A_{fw}), by optimizing m -step-ahead prediction of neural-behavioral data. This interconnected two-RNN system is visualized for $RNN1$ and $RNN1_{fw}$ in this computation graph. The superscript $\cdot^{(1)}$ indicating that parameters are for $RNN1$ and $RNN1_{fw}$ is omitted for simplicity.

a $\cdot^{(1)}$ or $\cdot^{(2)}$ superscript, e.g., $A^{(1)}$ and $A^{(2)}$. We provide the full two-section formulation for the model in appendix A.1.1 and the optimization details in appendix A.1.2. Briefly, each of the two learning stages consist of 2 optimizations, as follows:

Stage 1: Learning $RNN1$ and $RNN1_{fw}$

- 1a Learn $A^{(1)}$, $A_{fw}^{(1)}$, $K^{(1)}$, $K_{fw}^{(1)}$, and $C_z^{(1)}$, and extract latent states $x_k^{(1)}$ and $x_{k+m|k}^{(1)}$ (for $m > 1$) by minimizing the behavior prediction loss L_z from equation 4.
- 1b Learn $C_y^{(1)}$ by predicting neural data from $x_k^{(1)}$ and $x_{k+m|k}^{(1)}$, while minimizing the neural prediction loss L_y from equation 4.

Stage 2: Learning $RNN2$ and $RNN2_{fw}$

- 2a Learn $A^{(2)}$, $A_{fw}^{(2)}$, $K^{(2)}$, $K_{fw}^{(2)}$, and $C_y^{(2)}$, and extract latent states $x_k^{(2)}$ and $x_{k+m|k}^{(2)}$ (for $m > 1$) by minimizing the neural loss L_y , while including outputs of stage 1b as part of the predictions.
- 2b Learn $C_z^{(2)}$ by predicting behavior from $x_k^{(2)}$ and $x_{k+m|k}^{(2)}$, while minimizing the behavior loss L_z from equation 4, and including outputs of stage 1a as part of the predictions.

The explicit dissociation of the relevant dynamics and the above two-stage optimization allow us to first preferentially learn the (low-dimensional) shared dynamics between the two observations i.e., the behaviorally relevant neural dynamics, in stage 1. Then in stage 2, we learn any residual neural dynamics. Note that this residual neural dynamics ($x_k^{(2)}$) depend on the behaviorally relevant dynamic ($x_k^{(1)}$) as depicted in figure 1a (see appendix A.1 for details). The optional stage 2 is of interest for explaining neural dynamics beyond the ones related to behavior. This multi-stage approach has similarities to Sani et al. (2024), but here we have: 1) additional signals (u_k), 2) different losses, 3) a forecasting RNN within each model section (figure 1b), and additional steps that are discussed in the next section.

3.3 DISSOCIATION OF BEHAVIOR-SPECIFIC DYNAMICS

Optimizing behavior prediction (stage 1a) given neural activity y_k and input u_k can lead to learning behavior dynamics that are predictable from the input but are not encoded in the recorded neural activity. Although learning such behavior-specific dynamics enhances behavior decoding, it poses an interpretation challenge for neuroscience applications because one would not know what part of the learned dynamics are represented in the recorded brain regions. As shown in Vahidi et al.

(2024) this may lead to a misinterpretation of input-driven behavior-specific dynamics as intrinsic dynamics of the recorded brain region. To mitigate this possibility, we develop two additional steps in our method, that can 1) exclude such behavior-specific dynamics from $\mathbf{x}_k^{(1)}$, and 2) learn them separately as a distinct latent states $\mathbf{x}_k^{(3)}$ (figure 1a). Details are provided in appendix A.1.3. Briefly, *first*, to exclude behavior-specific dynamics, we introduce an optional preprocessing stage that predicts behavior from neural data, and passes this neurally-predicted behavior to be used in stages 1a and 2b. This preprocessing step ensures that the behaviorally relevant states learned in stage 1 ($\mathbf{x}_k^{(1)}$) are encoded in recorded neural activity \mathbf{y}_k , which can be crucial for interpretability in neuroscience studies. In our analyses of the real datasets (section 4.2), we always include this preprocessing step. *Second*, to still be able to learn behavior-specific dynamics, we add an optional post-hoc learning step (i.e., stage 3) that fits $RNN3$ and $RNN3_{fw}$ to any unexplained behavior and learns these input-driven behavior-specific dynamics as a distinct latent state $\mathbf{x}_k^{(3)}$. As we show in simulations (see section A.5.1 and figure A.2), the preprocessing step can exclude non-encoded behavior dynamics; and when desired in an application, the optional stage 3 can learn such dynamics to offset any behavior decoding loss incurred due to the preprocessing step, while still maintaining the interpretability of the model. We did not apply this post-hoc step in our real data analyses (section 4.2).

3.4 INFERENCE AND EVALUATION METRICS

After learning parameters of BRAID, we can readily use the learned mappings \mathbf{A} , \mathbf{K} (and \mathbf{A}_{fw} and \mathbf{K}_{fw}) to infer the 1-(and multi)-step-ahead predicted states \mathbf{x}_k (and $\mathbf{x}_{k+m|k}$) using equations 2 (and 3) for the held-out test data. Predicted neural activity and behavior are obtained by applying their corresponding decoders \mathbf{C}_y and \mathbf{C}_z to these inferred states. We also use the term “decoding” for behavior predictions because our model predicts behavior only using neural data and inputs, and never using behavior itself. To evaluate the performance of our models, we perform 5- and 2-fold cross-validation, for real data and simulation analyses, respectively. We calculate and report the Pearson Correlation Coefficient (CC) and in some cases (table A.8) also the coefficient of determination (R^2) between the predicted and actual observation, averaged across the data dimensions. We further report the m -step-ahead prediction accuracy, which gives a measure of how well the intrinsic dynamics \mathbf{A}_{fw} are learned. For simulation analyses with linear recursions (\mathbf{A}_{fw}), we additionally evaluate the learned intrinsic dynamics by comparing the eigenvalues of \mathbf{A}_{fw} between the true and learned model (see appendix A.1.5).

4 EXPERIMENTAL RESULTS

4.1 SIMULATION EXPERIMENTS

We validated BRAID in three simulations with different nonlinear neural-behavioral-input structures, to show that it can learn intrinsic behaviorally relevant neural dynamics in presence of inputs.

4.1.1 BRAID ACHIEVES NEAR OPTIMAL NEURAL-BEHAVIORAL PREDICTIVE ACCURACY IN NONLINEAR INPUT-DRIVEN SIMULATIONS

First, we considered an input-driven dynamical system as in equation 1, but with only the behavior mapping \mathbf{C}_z being nonlinear with the mapping $f_{C_z}(\nu) := a \sin(\nu) + b\nu$, as detailed in section A.4.2 (figure 2a). We generated 10 random parameter sets as our true models and generated data from them. First, we implemented an automatic selection of nonlinearity for BRAID by setting each of \mathbf{A} , \mathbf{K} , \mathbf{C}_y , or \mathbf{C}_z to linear or nonlinear, resulting in 2^4 different BRAID models, and finding the model with the best behavior decoding in the training data. Across all 10 realizations and 2 cross-validated folds, the behavior decoder, setting \mathbf{C}_z to be nonlinear was correctly identified as the best performing nonlinearity in 100% of the cases. Additionally, we evaluated BRAID with nonlinearity only in one of \mathbf{A} , \mathbf{K} , \mathbf{C}_y , or \mathbf{C}_z . The model with nonlinear behavior decoder \mathbf{C}_z outperformed other nonlinearity choices as well as linear models i.e., IPSID and the fully linear BRAID, in behavior decoding and neural prediction. BRAID further outperformed DPAD, which is nonlinear but does not account for the input u_k , in neural-behavioral prediction. In fact, both BRAID with nonlinear \mathbf{C}_z and BRAID with automatic nonlinearity selection achieved almost the

same neural-behavioral prediction as the true simulated models, demonstrating BRAID’s success in accurately learning the nonlinear input-driven dynamical system (table 2).

Table 2: 1-step-ahead prediction results for nonlinear simulation with sinusoidal behavior mapping. Mean \pm s.e.m. is across 20 runs (10 datasets, 2 folds). State dimension is always set to ground truth. True model’s outcome indicates the “Ideal” accuracy. **Bold**: within 1 s.e.m. of ideal.

Method	Behavior decoding CC	Neural prediction CC
IPSID	0.4567 \pm 0.0527	0.8901 \pm 0.0304
linear BRAID	0.4558 \pm 0.0528	0.8893 \pm 0.0304
DPAD <i>Nonlin C_z</i>	0.4958 \pm 0.0620	0.3767 \pm 0.0680
BRAID <i>Nonlin A</i>	0.4735 \pm 0.0572	0.8127 \pm 0.0528
BRAID <i>Nonlin K</i>	0.7680 \pm 0.0467	0.6983 \pm 0.0473
BRAID <i>Nonlin C_y</i>	0.4558 \pm 0.0528	0.8887 \pm 0.0305
BRAID <i>Nonlin C_z</i>	0.8696 \pm 0.0487	0.8913 \pm 0.0305
BRAID <i>Auto Nonlin</i>	0.8693 \pm 0.0487	0.8913 \pm 0.0305
True model (ideal)	0.8737 \pm 0.0486	0.8921 \pm 0.0306

4.1.2 BRAID DISSOCIATES INTRINSIC DYNAMICS FROM INPUT DYNAMICS IN SIMULATIONS

Next, we sought to validate BRAID’s ability to disentangle intrinsic and input-driven contributions to neural dynamics. BRAID simultaneously learns a predictor form and a generative form representation of the dynamics, the latter of which directly describes the intrinsic dynamics in terms of the mapping $\mathbf{A}_{fw}(\cdot)$ (section 3.1, figure 1b). In this simulation, we kept the ground truth state transitions linear so that we could precisely quantify the intrinsic dynamics and their learning error via the eigenvalues of the state transition matrix \mathbf{A}_{fw} (appendix A.1.5). We analyzed data from three sets of simulated dynamical systems with distinct nonlinear structures (see appendix A.4 for details): (1) Spiral behavior manifold (figure 2a-d), (2) trigonometric behavior manifold (also explained in section 4.1.1, figure 2e-h), (3) trigonometric input-encoder (figure A.1). For each simulation, we generated realizations from 10 different systems with randomly generated sets of parameters. Across all three simulations, BRAID accurately learned the intrinsic dynamics, resulting in smaller error in eigenvalues of the transition matrix, compared with DPAD, which is nonlinear but does not consider the input, and compared with linear BRAID and IPSID, which consider input, but are linear (figures 2c,g, A.1d). We demonstrate in an ablation study that learning a separate generative model for forecasting is crucial for correct learning of the intrinsic dynamics (table A.3). This more accurate intrinsic dynamics coupled with the input also resulted in BRAID achieving better behavior decoding (figures 2b,f, A.1a) as well as neural prediction (figure A.1c) compared to baselines. These results suggest that failing to account for either nonlinearity or input may lead to less accurate models and a misinterpretation of the intrinsic dynamics in nonlinear data.

Another metric for how well intrinsic dynamics are learned is forecasting, where behavior is predicted multiple steps into the future, without observing new neural data and only by observing the future input (section 3.1). Forecasting evolves the state dynamics according to the learned intrinsic dynamics (equation 3) and as such validates their accurate learning. We performed forecasting up to 32 steps ahead and found that the nonlinear models with input consistently outperformed linear models as well as the nonlinear DPAD, which does not consider input (figures 2d,h and A.1e,f).

4.2 NON-HUMAN PRIMATE MOTOR CORTICAL ACTIVITY DURING REACHING

We applied our method to a publicly available dataset recorded from a non-human primate (NHP) performing reaching movements O’Doherty et al. (2017) (figure 3a). We took either the smoothed spike counts or raw LFP from primary motor cortex (M1) as neural time-series \mathbf{y}_k , fingertip’s position and velocity as behavior \mathbf{z}_k , and sensory task instructions (target location) as the input \mathbf{u}_k (appendix A.3). Sensory inputs can have their own dynamics, which are distinct from the intrinsic dynamics of the motor cortex. Our goal is to learn the intrinsic dynamics in M1 related to movement while disentangling them from the dynamics of sensory input and also from any behavior-specific dynamics. Therefore, we importantly include BRAID’s behavior preprocessing stage (section 3.3).

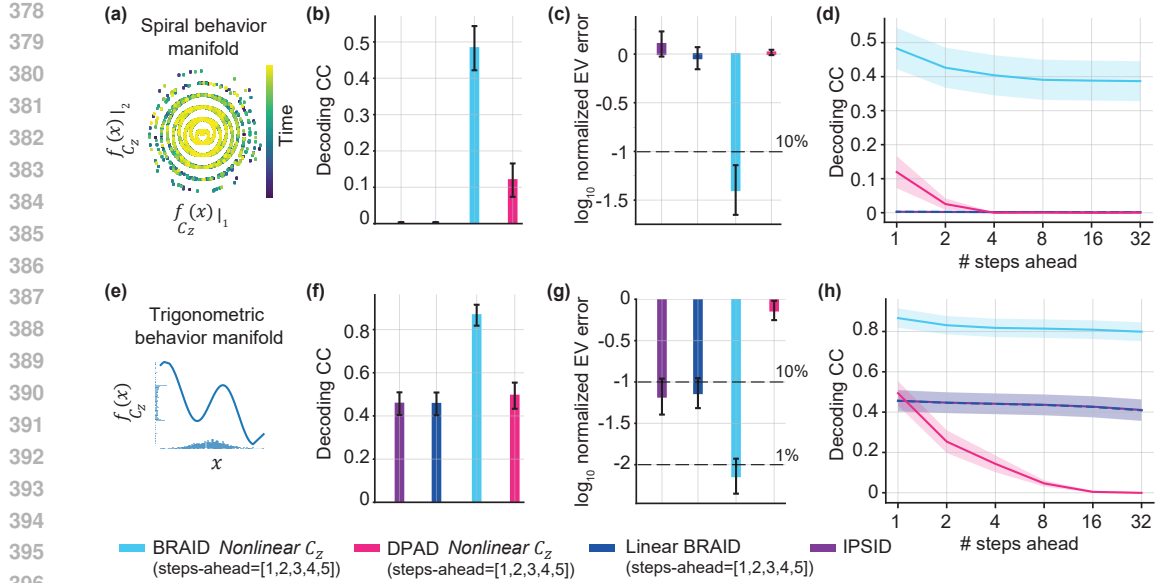


Figure 2: **BRAID better learns the intrinsic shared dynamics by simultaneously modeling input and nonlinearity, and by optimizing forecasting.** (a-d) Results for simulation with spiral behavioral manifold. (b) 1-step-ahead behavior decoding for nonlinear BRAID, nonlinear DPAD, linear BRAID, and IPSID (c). Error in identifying intrinsic dynamics of the true model quantified by the error in learning the eigenvalues of \mathbf{A}_{fw} . (d) Behavior decoding forecasts for 1 to 32 steps ahead, enabled by learning the intrinsic dynamics (\mathbf{A}_{fw}), with predictions optimized for [1, 2, 3, 4, 5]-steps-ahead (section 3). (e-h) Same as (a-d) for simulation with trigonometric behavior mapping.

We fitted BRAID with different nonlinearity choices as in our simulations: (1) nonlinear recursion $\mathbf{A}(\cdot)/\mathbf{A}_{fw}(\cdot)$, (2) nonlinear encoder $\mathbf{K}(\cdot)/\mathbf{K}_{fw}(\cdot)$, (3) nonlinear decoders $\mathbf{C}_y(\cdot)$ and $\mathbf{C}_z(\cdot)$, and a fully linear variant, linear BRAID. We included $m = [1, 2, 4, 8]$ -steps-ahead predictions in the BRAID loss, which for $m > 1$ engage the intrinsic behaviorally relevant dynamics (\mathbf{A}_{fw}) and allow their learning. To evaluate the learned intrinsic dynamics, we report neural-behavioral forecasting accuracy for different step-ahead horizons (figures 3b-c, and A.3 for LFP), while tabulating the 4-steps-ahead (i.e., 200ms) results to highlight BRAID’s advantage in forecasting (tables 3 and A.4). In addition to these results in the low-dimensional regime (stage 1 only, $n_x = n_1 = 16$), we also report results in the high-dimensional regime (both stages, $n_x = 64$, $n_1 = 16$) (table 3, figures A.5 and A.6). Among nonlinearity configurations, BRAID with nonlinear decoders provided the best fit to neural-behavioral data (table A.4). Moreover, BRAID’s behavior forecasting performance improved as more neurons were included (table A.7).

Next, we compared BRAID’s neural-behavioral forecasting to several ablation baselines (table 3 and figures 3 and A.5). First, BRAID outperformed linear BRAID, i.e., a similar but fully linear model. Second, BRAID outperformed DPAD (Sani et al., 2024), which can have decoder nonlinearities but does not consider inputs. BRAID’s advantage shows the importance of considering the effects of sensory inputs on neural-behavioral dynamics. Third, we compared to U-BRAID, which removes the first stage of BRAID and thus loses prioritization. BRAID consistently outperformed U-BRAID in behavioral forecasting, but did not match the neural forecasting of U-BRAID unless it was given enough latent state dimensions (table 3, $n_x = 64$), which is expected given U-BRAID’s singular objective being neural prediction. This comparison shows the benefit of prioritization for learning low-dimensional representations of intrinsic behaviorally relevant dynamics, and confirms that BRAID’s stage 2, BRAID can capture any remaining non-behavioral neural dynamics. Finally, BRAID’s low dimensional latent state trajectories were better separated for different movement directions compared to those of U-BRAID and DPAD, suggesting that BRAID’s latent states are more congruent with behavior, i.e., more behaviorally relevant (figure A.4).

We also compared BRAID’s neural-behavioral forecasting performance to mmPLRNN, which models multi-modal data using piecewise-linear RNNs, and has the option to model inputs although prior

Table 3: Forecasting performance (4-step-ahead) compared to baselines in NHP dataset for models with low ($n_x = 16$) and high-dimensional ($n_x = 64$) latent states. $n_1 = 16$ for BRAID, linear BRAID, and DPAD. R^2 results were similar (table A.8). Tables A.4, A.5 and A.6 have additional results.

Method	Behavior forecasting CC		Neural forecasting CC	
	$n_x = 16$	$n_x = 64$	$n_x = 16$	$n_x = 64$
linear BRAID	0.7453 ± 0.0066	0.7409 ± 0.0059	0.1767 ± 0.0054	0.3784 ± 0.0078
DPAD	0.6706 ± 0.0096	0.7352 ± 0.0079	0.2067 ± 0.0062	0.3611 ± 0.0080
U-BRAID	0.7663 ± 0.0069	0.8049 ± 0.0068	0.4089 ± 0.0076	0.4185 ± 0.0074
mmPLRNN	0.6851 ± 0.0143	0.7328 ± 0.00361	0.3162 ± 0.0107	0.3570 ± 0.0223
BRAID (ours)	0.8042 ± 0.0085	0.7970 ± 0.0086	0.3274 ± 0.0078	0.4123 ± 0.0077

work had not explored this input option. BRAID outperformed input-driven mmPLRNN networks in forecasting both behavior and neural activity, across all forecasting horizons, with the exception of 2-steps-ahead neural prediction (table 3, figure 3). Note that BRAID’s better decoding is achieved despite the fact that mmPLRNN, by design, incorporates behavior *as an input* during inference, whereas BRAID does not. Additionally, we compared BRAID to TNDM (Hurwitz et al., 2021), a nonlinear sequential autoencoder that models neural-behavioral dynamics, but does not include the effect of external inputs. We extended TNDM beyond the original work to create a version that also adds sensory inputs as an additional input besides neural activity (appendix A.2.6). We analyzed non-smoothed spike counts from the same dataset and found that BRAID significantly outperformed TNDM in behavior decoding while achieving comparable neural prediction (table A.5). Extending TNDM to include inputs significantly improved its behavior decoding, but it still did not reach that of BRAID (table A.5). Finally, we compared BRAID with CEBRA (Schneider et al., 2023), which is a non-dynamic convolutional encoder with a contrastive loss on behavior (section A.2.8), and with LFADS (Pandarinath et al., 2018), which is an unsupervised sequential autoencoder (section A.2.7). In both cases, BRAID outperformed these baselines in neural-behavioral prediction (table A.6).

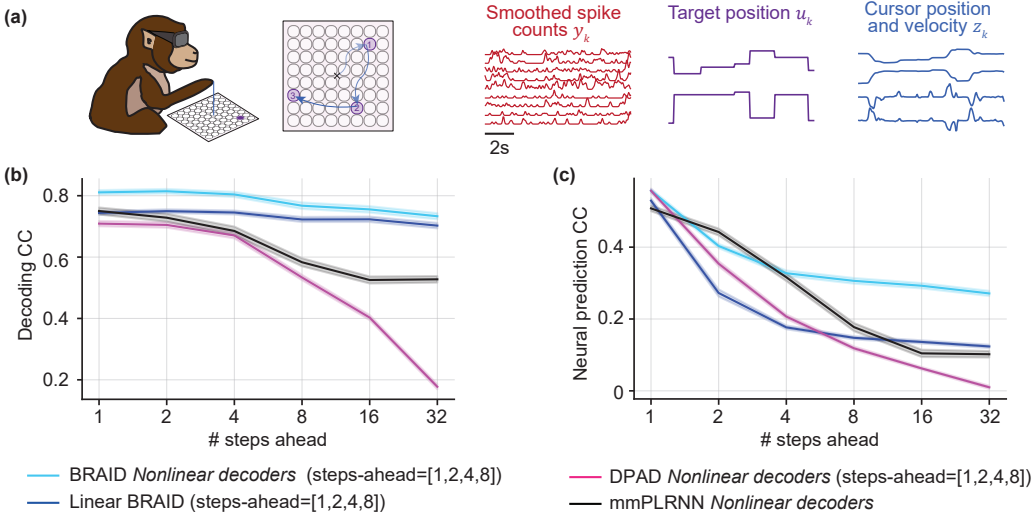


Figure 3: **BRAID outperforms baselines in neural-behavioral forecasting.** (a) Dataset and task visualization. (b) Behavior and, (c) neural activity forecasting correlation coefficient (CC) for BRAID, linear BRAID, DPAD, and mmPLRNN at low state dimension regime ($n_x=16$). Shaded areas show the s.e.m., across the 7 recording sessions and 5 cross-validation folds.

5 DISCUSSION

We introduced BRAID, a method for input-driven nonlinear dynamical modeling that disentangles intrinsic shared dynamics between two observation modalities from the effect of input. Here, we assume some external inputs are measured (e.g., from sensory stimuli or other brain regions) and

486 are available for modeling. This approach is distinct from the input-inference approach, where
487 unmeasured inputs are inferred from measured neural activity (Pandarinath et al., 2018; Schimel
488 et al., 2022). These two approaches are in a sense complementary. In our approach, any measured
489 inputs can be explicitly incorporated into the model to dissociate their dynamics from the intrinsic
490 dynamics of the measured neural-behavioral data. Practically, one cannot measure all inputs to a
491 given brain area, so our approach does not rule out the influence of unmeasured inputs on the learned
492 intrinsic dynamics. One could thus use the input-inference approach to infer such unmeasured inputs
493 from all measured signals. Note however that inferred inputs are ultimately a function of measured
494 signals and thus do not add any new information (unlike measured inputs), rather they can be thought
495 of as a decomposition of the measured signals based on certain assumptions (e.g. smoothness).

496 While the learning in BRAID is supervised by behavior, we only use neural activity and input (but
497 *not* behavior) during inference. This supervision allows stage 1 to extract behaviorally relevant
498 intrinsic dynamics with priority, while later stages learn neural-specific or behavior-specific dynam-
499 ics in independent optimizations. This multi-stage learning approach has similarities to some prior
500 works (Vahidi et al., 2024; Sani et al., 2024), but is fundamentally different from other works that
501 use a single multi-modal optimization loss (Kramer et al., 2022; Gondur et al., 2024; Hurwitz et al.,
502 2021), which may miss prioritization of the shared dynamics over unshared dynamics (Sani et al.,
503 2024). In fact, some of these works are focused on multi-modal inference and aim to fuse all shared
504 and unshared information into the same latent space (Kramer et al., 2022; Gondur et al., 2024). The
505 multi-stage approach avoids this fusion by focusing on shared dynamics during a dedicated first
506 optimization stage with only cross-modality prediction (e.g., behavior decoding) as the objective.

507 To disentangle input dynamics from intrinsic dynamic, BRAID consists of three stages, each learn-
508 ing a predictor and a generator model. In each stage, BRAID’s predictor model in that stage (e.g.,
509 *RNNI*) infers the latent states, which are subsequently employed as the input to compute m -step-
510 ahead predictions via the associated generative model (e.g., *RNNI_{fw}*). BRAID’s predictor and gener-
511 ative models are learned jointly to maximize the m -step-ahead log-likelihood. This has analogies
512 to the encoder-decoder architectures such as those commonly used in variational inference (Kingma
513 & Welling, 2013). Specifically, the predictor and generator RNNs in BRAID have roles similar to
514 those of the encoder and decoder in variational inference, respectively. However, in variational infer-
515 ence, the posterior distribution of the unobserved variables given the data is parametrized and a
516 part of the optimization loss aims to enforce that distribution on the inferred latent variables (Chung
517 et al., 2015; Krishnan et al., 2015; Fraccaro et al., 2016; Luk et al., 2024). In contrast, we do not
518 impose such parametrization on the latent states in BRAID. Developing variational methods with the
519 same multi-section architecture and multi-stage learning as BRAID is an interesting future direction.

519 Similar to many prior works (Chung et al., 2015; Krishnan et al., 2015; Fraccaro et al., 2016; Luk
520 et al., 2024), BRAID has a causal formulation and can perform inference by recursively inferring the
521 next sample of the latent state after each new observation is measured. This is distinct from some
522 nonlinear approaches in neuroscience (Gao et al., 2016; Pandarinath et al., 2018; Hernandez et al.,
523 2020; Hurwitz et al., 2021; Keshtkaran et al., 2022; Gondur et al., 2024; Karniol-Tambour et al.,
524 2024) that perform inference non-causally in time. As such BRAID may also be a good candidate
525 for real-time decoding applications such as brain-machine-interfaces.

526 BRAID, as presented, is designed for single-session settings. Extending it for cross-session gen-
527 eralization is an interesting future direction that can follow approaches used in the literature, such
528 as aligning neural manifolds across sessions (Farshchian et al., 2019), or incorporating session-
529 specific read-in and read-out mappings while keeping shared model parameters fixed (Pandarinath
530 et al., 2018). The latter allows training across multiple sessions and adapting to new sessions with
531 minimal data by learning session-specific matrices. Furthermore, BRAID’s ability to disentangle
532 intrinsic dynamics may also facilitate generalization across tasks with different sensory instructions.

533 A fundamental challenge for nonlinear latent state models is the fact that many alternative models
534 may explain the data equally well. As such, to evaluate the learned dynamics, we are limited to
535 computable quantities related to measured signals, most importantly m -step-ahead prediction of
536 neural and behavioral signals. While our results suggest that in our dataset having nonlinear decoders
537 provides superior performance compared to other nonlinearities, this might not be case in another
538 dataset. In practice, one can search for the optimal nonlinearity based on the desired metric as done
539 by BRAID. In the special case of linear models, all correct latent models share the same eigenvalues
(Katayama, 2006), enabling a direct evaluation of learned latent dynamics (figure 2).

6 REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we are sharing the code for BRAID along with a Python notebook demonstrating its usage. At the moment, we provide these in a temporary double blind repository at <https://anonymous.4open.science/r/BRAID-5ABF>, but we will provide them on Github upon publication. We also provide model architecture and training details in appendix A.1.4. Finally, the dataset we used (O’Doherty et al., 2017) is publicly available for anyone interested in reproducing the results reported in section 4.2.

REFERENCES

- Hamidreza Abbaspourazad, Mahdi Choudhury, Yan T Wong, Bijan Pesaran, and Maryam M Shanechi. Multiscale low-dimensional motor cortical state dynamics predict naturalistic reach-and-grasp behavior. *Nature communications*, 12(1):607, 2021.
- Hamidreza Abbaspourazad, Eray Erturk, Bijan Pesaran, and Maryam M Shanechi. Dynamical flexible inference of nonlinear latent factors and structures in neural population activity. *Nature Biomedical Engineering*, 8(1):85–108, 2024.
- Mehdi Aghagolzadeh and Wilson Truccolo. Inference and decoding of motor cortex low-dimensional dynamics via latent state-space models. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 24(2):272–282, 2015.
- William E Allen, Michael Z Chen, Nandini Pichamoorthy, Rebecca H Tien, Marius Pachitariu, Liqun Luo, and Karl Deisseroth. Thirst regulates motivated behavior through modulation of brainwide neural population dynamics. *Science*, 364(6437):eaav3932, 2019.
- Dean V Buonomano and Wolfgang Maass. State-dependent computations: spatiotemporal processing in cortical networks. *Nature Reviews Neuroscience*, 10(2):113–125, 2009.
- Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. *Advances in neural information processing systems*, 28, 2015.
- Mark M Churchland, John P Cunningham, Matthew T Kaufman, Justin D Foster, Paul Nuyujukian, Stephen I Ryu, and Krishna V Shenoy. Neural population dynamics during reaching. *Nature*, 487(7405):51–56, 2012.
- John P Cunningham and Byron M Yu. Dimensionality reduction for large-scale neural recordings. *Nature neuroscience*, 17(11):1500–1509, 2014.
- Daniel Durstewitz. A state space approach for piecewise-linear recurrent neural networks for identifying computational dynamics from neural measurements. *PLoS computational biology*, 13(6):e1005542, 2017.
- Daniel Durstewitz, Georgia Koppe, and Max Ingo Thurm. Reconstructing computational system dynamics from neural data with recurrent neural networks. *Nature Reviews Neuroscience*, 24(11):693–710, 2023.
- Tatiana A Engel and Nicholas A Steinmetz. New perspectives on dimensionality and variability from large-scale cortical dynamics. *Current opinion in neurobiology*, 58:181–190, 2019.
- Ali Farshchian, Juan A. Gallego, Joseph P. Cohen, Yoshua Bengio, Lee E. Miller, and Sara A. Solla. ADVERSARIAL DOMAIN ADAPTATION FOR STABLE BRAIN-MACHINE INTERFACES. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Hyx6Bi0qYm>.
- Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models with stochastic layers. *Advances in neural information processing systems*, 29, 2016.
- Yuanjun Gao, Evan W Archer, Liam Paninski, and John P Cunningham. Linear dynamical neural population models through nonlinear embeddings. *Advances in neural information processing systems*, 29, 2016.

- 594 Rabia Gondur, Usama Bin Sikandar, Evan Schaffer, Mikio Christian Aoi, and Stephen L Keeley.
595 Multi-modal gaussian process variational autoencoders for neural and behavioral data. In *Inter-*
596 *national Conference on Learning Representations*, 2024.
- 597 Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of*
598 *statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- 600 Daniel Hernandez, Antonio Khalil Moretti, Ziqiang Wei, Shreya Saxena, John Cunningham, and
601 Liam Paninski. Nonlinear evolution via spatially-dependent linear dynamics for electrophysiol-
602 ogy and calcium data, 2020. URL <https://arxiv.org/abs/1811.02459>.
- 603 Cole Hurwitz, Akash Srivastava, Kai Xu, Justin Jude, Matthew Perich, Lee Miller, and Matthias
604 Hennig. Targeted neural dynamical modeling. *Advances in Neural Information Processing Sys-*
605 *tems*, 34:29379–29392, 2021.
- 607 James J Jun, Nicholas A Steinmetz, Joshua H Siegle, Daniel J Denman, Marius Bauza, Brian Barbar-
608 its, Albert K Lee, Costas A Anastassiou, Alexandru Andrei, Çağatay Aydın, et al. Fully integrated
609 silicon probes for high-density recording of neural activity. *Nature*, 551(7679):232–236, 2017.
- 610 Jonathan C Kao, Paul Nuyujukian, Stephen I Ryu, Mark M Churchland, John P Cunningham, and
611 Krishna V Shenoy. Single-trial dynamics of motor cortex and their applications to brain-machine
612 interfaces. *Nature communications*, 6(1):7759, 2015.
- 613 Orren Karniol-Tambour, David M Zoltowski, E Mika Diamanti, Lucas Pinto, Carlos D Brody,
614 David W Tank, and Jonathan W Pillow. Modeling state-dependent communication between brain
615 regions with switching nonlinear dynamical systems. In *The Twelfth International Conference on*
616 *Learning Representations*, 2024.
- 618 Tohru Katayama. *Subspace Methods for System Identification*. Springer Science & Business Me-
619 dia, 2006. ISBN 978-1-84628-158-7. URL [https://link.springer.com/book/10.](https://link.springer.com/book/10.1007%2F1-84628-158-X)
620 [1007%2F1-84628-158-X](https://link.springer.com/book/10.1007%2F1-84628-158-X).
- 621 Mohammad Reza Keshtkaran, Andrew R Sedler, Raeed H Chowdhury, Raghav Tandon, Diya Bas-
622 rai, Sarah L Nguyen, Hansem Sohn, Mehrdad Jazayeri, Lee E Miller, and Chethan Pandarinath.
623 A large-scale neural network training framework for generalized estimation of single-trial popu-
624 lation dynamics. *Nature Methods*, 19(12):1572–1577, 2022.
- 625 Timothy D Kim, Thomas Z Luo, Jonathan W Pillow, and Carlos D Brody. Inferring latent dy-
626 namics underlying neural population activity via neural differential equations. In *International*
627 *Conference on Machine Learning*, pp. 5551–5561. PMLR, 2021.
- 629 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. 2017. URL
630 <http://arxiv.org/abs/1412.6980>.
- 631 Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013. URL <https://arxiv.org/abs/1312.6114>.
- 632 Dmitry Kobak, Wieland Brendel, Christos Constantinidis, Claudia E Feierstein, Adam Kepecs,
633 Zachary F Mainen, Xue-Lian Qi, Ranulfo Romo, Naoshige Uchida, and Christian K Machens.
634 Demixed principal component analysis of neural population data. *elife*, 5:e10989, 2016.
- 635 Daniel Kramer, Philine L Bommer, Carlo Tombolini, Georgia Koppe, and Daniel Durstewitz. Re-
636 constructing nonlinear dynamical systems from multi-modal time series. In *Proceedings of the*
637 *39th International Conference on Machine Learning*, volume 162, pp. 11613–11633. PMLR, 17–
638 23 Jul 2022.
- 639 Rahul G. Krishnan, Uri Shalit, and David Sontag. Deep kalman filters, 2015. URL [https://](https://arxiv.org/abs/1511.05121)
640 arxiv.org/abs/1511.05121.
- 641 Enoch Luk, Eviatar Bach, Ricardo Baptista, and Andrew Stuart. Learning optimal filters using
642 variational inference, 2024. URL <https://arxiv.org/abs/2406.18066>.
- 643 Valerio Mante, David Sussillo, Krishna V Shenoy, and William T Newsome. Context-dependent
644 computation by recurrent dynamics in prefrontal cortex. *nature*, 503(7474):78–84, 2013.

- 648 Erfan Nozari, Maxwell A Bertolero, Jennifer Stiso, Lorenzo Caciagli, Eli J Cornblath, Xiaosong
649 He, Arun S Mahadevan, George J Pappas, and Dani S Bassett. Macroscopic resting-state brain
650 dynamics are best described by linear models. *Nature biomedical engineering*, 8(1):68–84, 2024.
651
- 652 Joseph E. O’Doherty, Mariana M. B. Cardoso, Joseph G. Makin, and Philip N. Sabes. Nonhuman
653 Primate Reaching with Multichannel Sensorimotor Cortex Electrophysiology, May 2017. URL
654 <https://doi.org/10.5281/zenodo.583331>.
- 655 Dan O’Shea and Chethan Pandarinath. Lfads run manager, December 2021. URL <https://doi.org/10.5281/zenodo.5790161>.
656
657
- 658 Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky,
659 Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg,
660 et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature*
661 *methods*, 15(10):805–815, 2018.
- 662 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
663 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
664 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
665 12:2825–2830, 2011.
666
- 667 Evan D Remington, Seth W Egger, Devika Narain, Jing Wang, and Mehrdad Jazayeri. A dynamical
668 systems perspective on flexible motor timing. *Trends in cognitive sciences*, 22(10):938–952, 2018.
- 669 Virginia Rutten, Alberto Bernacchia, Maneesh Sahani, and Guillaume Hennequin. Non-reversible
670 gaussian processes for identifying latent dynamical structure in neural data. *Advances in neural*
671 *information processing systems*, 33:9622–9632, 2020.
672
- 673 Omid G Sani, Hamidreza Abbaspourazad, Yan T Wong, Bijan Pesaran, and Maryam M Shanechi.
674 Modeling behaviorally relevant neural dynamics enabled by preferential subspace identification.
675 *Nature Neuroscience*, 24(1):140–149, 2021.
- 676 Omid G Sani, Bijan Pesaran, and Maryam M Shanechi. Dissociative and prioritized modeling of
677 behaviorally relevant neural dynamics using recurrent neural networks. *Nature Neuroscience*, pp.
678 1–13, 2024.
679
- 680 Britton A Sauerbrei, Jian-Zhong Guo, Jeremy D Cohen, Matteo Mischiati, Wendy Guo, Mayank
681 Kabra, Nakul Verma, Brett Mensh, Kristin Branson, and Adam W Hantman. Cortical pattern
682 generation during dexterous movement is input-driven. *Nature*, 577(7790):386–391, 2020.
- 683 Marine Schimel, Ta-Chu Kao, Kristopher T Jensen, and Guillaume Hennequin. ilqr-vae: control-
684 based learning of input-driven dynamics with applications to neural data. biorxiv. In *International*
685 *Conference on Learning Representations*, 2022.
686
- 687 Steffen Schneider, Jin Hwa Lee, and Mackenzie Weygandt Mathis. Learnable latent embeddings for
688 joint behavioural and neural analysis. *Nature*, 617(7960):360–368, 2023.
- 689 Jeffrey S Seely, Matthew T Kaufman, Stephen I Ryu, Krishna V Shenoy, John P Cunningham,
690 and Mark M Churchland. Tensor analysis reveals distinct population structure that parallels the
691 different computational roles of areas m1 and v1. *PLoS computational biology*, 12(11):e1005164,
692 2016.
693
- 694 Krishna V Shenoy and Jonathan C Kao. Measurement, manipulation and modeling of brain-wide
695 neural population dynamics. *Nature communications*, 12(1):633, 2021.
- 696 Nicholas A Steinmetz, Peter Zátka-Haas, Matteo Carandini, and Kenneth D Harris. Distributed
697 coding of choice, action and engagement across the mouse brain. *Nature*, 576(7786):266–273,
698 2019.
699
- 700 Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Charu Bai Reddy, Matteo Carandini, and
701 Kenneth D Harris. Spontaneous behaviors drive multidimensional, brainwide activity. *Science*,
364(6437):eaav7893, 2019.

- 702 Damar Susilaradeya, Wei Xu, Thomas M Hall, Ferran Galán, Kai Alter, and Andrew Jackson. Ex-
703 trinsic and intrinsic dynamics in movement intermittency. *Elife*, 8:e40145, 2019.
704
- 705 David Sussillo, Sergey D Stavisky, Jonathan C Kao, Stephen I Ryu, and Krishna V Shenoy. Making
706 brain-machine interfaces robust to future neural variability. *Nature communications*, 7(1):13749,
707 2016.
- 708 Parsa Vahidi, Omid G Sani, and Maryam M Shanechi. Modeling and dissociation of intrinsic
709 and input-driven neural population dynamics underlying behavior. *Proceedings of the National
710 Academy of Sciences*, 121(7):e2212887121, 2024.
711
- 712 Peter Van Overschee and Bart De Moor. *Subspace Identification for Linear Systems*. Springer US,
713 Boston, MA, 1996. ISBN 978-1-4613-8061-0. doi: 10.1007/978-1-4613-0465-4.
- 714 Saurabh Vyas, Matthew D Golub, David Sussillo, and Krishna V Shenoy. Computation through
715 neural population dynamics. *Annual review of neuroscience*, 43(1):249–275, 2020.
716
- 717 Anqi Wu, Nicholas A Roy, Stephen Keeley, and Jonathan W Pillow. Gaussian process based non-
718 linear latent structure discovery in multivariate spike train data. *Advances in neural information
719 processing systems*, 30, 2017.
- 720 Yuxiao Yang, Omid G Sani, Edward F Chang, and Maryam M Shanechi. Dynamic network model-
721 ing and dimensionality reduction for human ecog activity. *Journal of neural engineering*, 16(5):
722 056014, 2019.
- 723 Yuxiao Yang, Shaoyu Qiao, Omid G Sani, J Isaac Sedillo, Breonna Ferrentino, Bijan Pesaran, and
724 Maryam M Shanechi. Modelling and prediction of the dynamic responses of large-scale brain net-
725 works during direct electrical stimulation. *Nature biomedical engineering*, 5(4):324–345, 2021.
726
- 727 Ding Zhou and Xue-Xin Wei. Learning identifiable and interpretable latent models of high-
728 dimensional neural activity using pi-vae. *Advances in Neural Information Processing Systems*,
729 33:7234–7247, 2020.

731 A APPENDIX

732 A.1 METHOD DETAILS

733 A.1.1 TWO-SECTION FORMULATION

734 In equations 1, 2, and 3, we combined both latent state sections of our model ($\mathbf{x}_k^{(1)}$ and $\mathbf{x}_k^{(2)}$) for
735 simpler exposition. Here, we present the complete two-section formulation. The predictor form part
736 of the model (equation 2) can be written as follows:

$$\begin{cases} \begin{bmatrix} \mathbf{x}_{k+1|k}^{(1)} \\ \mathbf{x}_{k+1|k}^{(2)} \end{bmatrix} &= \begin{bmatrix} \mathbf{A}^{(1)}(\mathbf{x}_{k|k-1}^{(1)}) \\ \mathbf{A}^{(2)}(\mathbf{x}_{k|k-1}^{(2)}) \end{bmatrix} + \begin{bmatrix} \mathbf{K}^{(1)}(\mathbf{y}_k, \mathbf{u}_k) \\ \mathbf{K}^{(2)}(\mathbf{y}_k, \mathbf{u}_k, \mathbf{x}_{k|k-1}^{(1)}) \end{bmatrix} \\ \hat{\mathbf{y}}_{k|k-1} &= \mathbf{C}_y^{(1)}(\mathbf{x}_{k|k-1}^{(1)}, \mathbf{u}_k) + \mathbf{C}_y^{(2)}(\mathbf{x}_{k|k-1}^{(2)}, \mathbf{u}_k) \\ \hat{\mathbf{z}}_{k|k-1} &= \mathbf{C}_z^{(1)}(\mathbf{x}_{k|k-1}^{(1)}, \mathbf{u}_k) + \mathbf{C}_z^{(2)}(\mathbf{x}_{k|k-1}^{(2)}, \mathbf{u}_k) \end{cases} \quad (\text{A.1})$$

737 where we have also included the two-section formulation for the prediction of observations as lines
738 2-3 of the equation. As before, $\mathbf{y}_k \in \mathbb{R}^{n_y}$ and $\mathbf{z}_k \in \mathbb{R}^{n_z}$ are the observed high-dimensional neural
739 activity and behavior respectively while $\mathbf{u}_k \in \mathbb{R}^{n_u}$ represents the measured inputs to the dynamical
740 system. Here, the overall latent state, $\mathbf{x}_k \in \mathbb{R}^{n_x}$, which describes the dynamics underlying the
741 neural-behavioral data, is constructed such that the behaviorally relevant neural dynamics, repre-
742 sented by $\mathbf{x}_k^{(1)} \in \mathbb{R}^{n_1}$, are dissociated from the the irrelevant ones, represented by $\mathbf{x}_k^{(2)} \in \mathbb{R}^{n_x - n_1}$.
743

744 The predictor form RNNs in equation A.1 (i.e., $RNN1$ and $RNN2$) are complemented by another
745 set of RNNs (i.e., $RNN1_{fw}$ and $RNN2_{fw}$) that constitute the generative form part of the model
(equation 3), and enable m -step-ahead (for $m > 1$) prediction of latent states and neural-behavioral

756 data. The generative RNNs were again shown with a combined latent state in equation 3 for simpler
 757 exposition. The following equations show the complete two-section formulation:
 758

$$\begin{cases} \begin{bmatrix} \mathbf{x}_{k+m|k}^{(1)} \\ \mathbf{x}_{k+m|k}^{(2)} \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_{fw}^{(1)}(\mathbf{x}_{k+m-1|k}^{(1)}) \\ \mathbf{A}_{fw}^{(2)}(\mathbf{x}_{k+m-1|k}^{(2)}) \end{bmatrix} + \begin{bmatrix} \mathbf{K}_{fw}^{(1)}(\mathbf{u}_{k+m-1}) \\ \mathbf{K}_{fw}^{(2)}(\mathbf{u}_{k+m-1}, \mathbf{x}_{k+m-1|k}^{(1)}) \end{bmatrix} \\ \hat{\mathbf{y}}_{k+m|k} &= \mathbf{C}_y^{(1)}(\mathbf{x}_{k+m|k}^{(1)}, \mathbf{u}_{k+m}) + \mathbf{C}_y^{(2)}(\mathbf{x}_{k+m|k}^{(2)}, \mathbf{u}_{k+m}) \\ \hat{\mathbf{z}}_{k+m|k} &= \mathbf{C}_z^{(1)}(\mathbf{x}_{k+m|k}^{(1)}, \mathbf{u}_{k+m}) + \mathbf{C}_z^{(2)}(\mathbf{x}_{k+m|k}^{(2)}, \mathbf{u}_{k+m}) \end{cases} \quad (\text{A.2})$$

766 where $m > 1$, and $\mathbf{x}_{k+1|k}^{(1)}$ and $\mathbf{x}_{k+1|k}^{(2)}$ (i.e., $m = 1$) are taken from equation A.1. A visualization
 767 of how the formulations in equations A.1 and A.2 are connected is provided in figure 1. In equa-
 768 tion A.2, we have also included the two-section formulation for the prediction of neural-behavioral
 769 observations as lines 2-3 of the equation, showing that applying the same decoders $\mathbf{C}_y^{(1)}/\mathbf{C}_y^{(2)}$ and
 770 $\mathbf{C}_z^{(1)}/\mathbf{C}_z^{(2)}$ as in equation A.1 to the m -step-ahead predicted latents $\mathbf{x}_{k+m|k}^{(1)}/\mathbf{x}_{k+m|k}^{(2)}$ gives the m -
 771 step-ahead predictions of the neural-behavioral data ($\hat{\mathbf{y}}_{k+m|k}$ and $\hat{\mathbf{z}}_{k+m|k}$).
 772
 773

774 Equations A.1 and A.2 together constitute the two-section formulation of the BRAID model, which
 775 consists of 12 transformations in total: $\mathbf{A}(\cdot)$, $\mathbf{A}_{fw}(\cdot)$, $\mathbf{K}(\cdot)$, $\mathbf{K}_{fw}(\cdot)$, $\mathbf{C}_z(\cdot)$, and $\mathbf{C}_y(\cdot)$, each having
 776 two sections denoted with the $\cdot^{(1)}$ and $\cdot^{(2)}$ superscripts.
 777

778 A.1.2 LEARNING ALGORITHM STEPS

779 In sections A.1.2-A.1.3, we provide detailed formulations of the optimization stages used in BRAID
 780 during learning. For simplicity, we explain the optimizations in terms of 1-step ahead predictions,
 781 which involve predictor form parameters of the model. Formulations for m -step-ahead predictions,
 782 which constitute additional terms in the overall loss (equations 4), are analogous to those provided
 783 here, but instead of the predictor form parameters (equation A.1) they engage the generative form
 784 parameters (equation A.2).
 785

786 Note that regardless of what step-ahead predictions were included during training, the learned model
 787 can be used to predict the latent state and neural-behavioral data at m -steps ahead for any desired
 788 m using equations 2 and 3 (or A.1 and A.2) together. For example, in figure 3, only [1, 2, 4, 8]
 789 step ahead predictions are included in the optimization loss (equation 4, but we evaluate the learned
 790 models with predictions up to 32-steps ahead.

791 We develop a two-stage optimization algorithm for learning parameters of the two sections of the
 792 BRAID model (equation A.1). We note that the following 2 stages are sequential. This means that
 793 parameters associated with behaviorally relevant states ($\mathbf{x}_k^{(1)}$) are fully learned with *RNN1*, then if
 794 needed, the remaining parameters corresponding to non-relevant dynamics ($\mathbf{x}_k^{(2)}$) can be learned via
 795 *RNN2*. In all the optimizations described below, we use the mean-squared-error (MSE) of predicting
 796 observations as the loss function, but we note that the MSE is proportional to the negative log-
 797 likelihood (NLL) for isotropic Gaussian-distributed data. Below we provide the details for the 4
 798 optimizations that are performed in the two learning stages of BRAID.
 799

800 Stage 1:

801 **1a** First, BRAID learns a recurrent neural network (*RNN1*) with n_1 states, to minimize behavior
 802 prediction MSE given past neural data and inputs (equation A.3). This ensures that *RNN1* only
 803 learns neural dynamics that are relevant to (i.e., predictive of) behavior. The states of *RNN1*
 804 constitute the first set of latent states in the BRAID model: $\mathbf{x}_k^{(1)}$. This optimization step can be
 805 formulated as:
 806

$$\begin{cases} \mathbf{x}_{k+1}^{(1)} &= \mathbf{A}^{(1)}(\mathbf{x}_k^{(1)}) + \mathbf{K}^{(1)}(\mathbf{y}_k, \mathbf{u}_k) \\ \mathbf{z}_k &= \mathbf{C}_z^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k) \\ \text{loss} &: \text{MSE}(\mathbf{z}_k, \mathbf{C}_z^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.3})$$

810 **1b** Next, in a second optimization, we learn a transformation $\mathbf{C}_y^{(1)}(\cdot)$ that maps $\mathbf{x}_k^{(1)}$ to neural activity
 811 while minimizing neural prediction MSE (equation A.4):
 812

$$813 \begin{cases} \mathbf{y}_k &= \mathbf{C}_y^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k) \\ 814 \text{loss} &: \text{MSE}(\mathbf{y}_k, \mathbf{C}_y^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.4})$$

816 The above 2 steps conclude stage 1 of learning, i.e., learning intrinsic behaviorally relevant
 817 dynamics $\mathbf{x}_k^{(1)}$. Next we explain the (optional) remaining stage 2, which can learn any remaining
 818 dynamics in neural activity \mathbf{x}_k^2 .
 819

820 **Stage 2:**

821 **2a** We learn a second recurrent neural network (*RNN2*) with $n_2 := n_x - n_1$ states, to minimize
 822 the MSE loss of predicting the residual neural activity, i.e., $\mathbf{y}'_k := \mathbf{y}_k - \mathbf{C}_y^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k)$, given
 823 past neural activity and inputs (equation A.5). States of *RNN2*, i.e., $\mathbf{x}_k^{(2)}$, together with $\mathbf{x}_k^{(1)}$ from
 824 stage 1 constitute the full neural dynamics, i.e., $\mathbf{x}_k = [\mathbf{x}_k^{(1)} \quad \mathbf{x}_k^{(2)}]^T$. This optimization step can
 825 be formulated as:
 826
 827

$$828 \begin{cases} \mathbf{x}_{k+1}^{(2)} &= \mathbf{A}^{(2)}(\mathbf{x}_k^{(2)}) + \mathbf{K}^{(2)}(\mathbf{y}_k, \mathbf{u}_k, \mathbf{x}_k^{(1)}) \\ 829 \mathbf{y}'_k &= \mathbf{C}_y^{(2)}(\mathbf{x}_k^{(2)}, \mathbf{u}_k) \\ 830 \text{loss} &: \text{MSE}(\mathbf{y}'_k, \mathbf{C}_y^{(2)}(\mathbf{x}_k^{(2)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.5})$$

833 **2b** Finally, another readout, $\mathbf{C}_z^{(2)}$, can be learned to map $\mathbf{x}_k^{(2)}$ to the residual behavior, i.e., $\mathbf{z}'_k :=$
 834 $\mathbf{z}_k - \mathbf{C}_z^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k)$, to minimizing the overall behavioral loss (equation A.6):
 835

$$836 \begin{cases} \mathbf{z}'_k &= \mathbf{C}_z^{(2)}(\mathbf{x}_k^{(2)}, \mathbf{u}_k) \\ 837 \text{loss} &: \text{MSE}(\mathbf{z}'_k, \mathbf{C}_z^{(2)}(\mathbf{x}_k^{(2)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.6})$$

839 Note that the optimization in stage 2b does not change *RNN2* or $\mathbf{x}_k^{(2)}$ that were learned in stage
 840 2a. So although stage 2b is supervised by behavior, the second set of states $\mathbf{x}_k^{(2)}$ are still learned
 841 unsupervised with respect to behavior.
 842

843 In case a very low state dimension is specified by the user for stage 1 (n_1 lower than the ground truth
 844 shared dimensionality), *RNN1* would not have enough capacity to learn all behaviorally relevant
 845 dynamics. In that case, some behaviorally relevant neural dynamics will be left for *RNN2* in stage
 846 2 to learn. This is why the $\mathbf{C}_z^{(2)}$ transformation from $\mathbf{x}_k^{(2)}$ to behavior is included in the model, to
 847 allow such behaviorally relevant information in $\mathbf{x}_k^{(2)}$ to be utilized to improve behavior decoding.
 848

849 A.1.3 NON-ENCODED BEHAVIOR-SPECIFIC DYNAMICS

850 To remove the non-encoded behavior-specific dynamics, as a preprocessing step, we fit a high-
 851 dimensional ($n_x = 150$ in all real data analyses) unsupervised RNN to extract neural dynamics
 852 alone by minimizing neural prediction MSE (equation A.7):
 853

$$854 \begin{cases} \mathbf{x}_{k+1}^{(0)} &= \mathbf{A}^{(0)}(\mathbf{x}_k^{(0)}) + \mathbf{K}^{(0)}(\mathbf{y}_k, \mathbf{u}_k) \\ 855 \mathbf{y}_k &= \mathbf{C}_y^{(0)}(\mathbf{x}_k^{(0)}, \mathbf{u}_k) \\ 856 \text{loss} &: \text{MSE}(\mathbf{y}_k, \mathbf{C}_y^{(0)}(\mathbf{x}_k^{(0)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.7})$$

859 Then a readout $\mathbf{C}_z^{(0)}$ is trained to map these neurally relevant states $\mathbf{x}_k^{(0)}$ to behavior:
 860

$$861 \begin{cases} \mathbf{z}_k &= \mathbf{C}_z^{(0)}(\mathbf{x}_k^{(0)}) \\ 862 \text{loss} &: \text{MSE}(\mathbf{z}_k, \mathbf{C}_z^{(0)}(\mathbf{x}_k^{(0)})) \end{cases} \quad (\text{A.8})$$

After parameters of the above are learned, we run inference on the training data to obtain the filtered behavior as output of the preprocessing RNN model (first line in equation A.8) and subsequently use it in place of the original behavior in BRAID (in equations A.1, A.3, A.6). In simulations, we validate that this additional stage can successfully remove any input-driven behavior dynamics not encoded in the neural recordings (figure A.2). We include this preprocessing step in all reported real data analyses with BRAID.

Stage 3: As mentioned in 3.3, if desired, BRAID can also learn behavior-specific dynamics as separate dissociated latent states using a post-hoc learning step. This step is performed after BRAID’s main learning is done and is meant to be used in conjunction with BRAID’s preprocessing stage explained above. In this post-hoc learning step (stage 3), we first infer the behavior using the originally learned BRAID model. We then obtain the residual behavior (\mathbf{z}''_k) by subtracting the inferred behavior from the measured behavior. We then learn a third RNN (*RNN3*) that optimizes the prediction of the residual behavior using *only* the external inputs. The following equations summarize this step:

$$\begin{cases} \mathbf{z}''_k & := \mathbf{z}_k - [\mathbf{C}_z^{(1)}(\mathbf{x}_k^{(1)}, \mathbf{u}_k) - \mathbf{C}_z^{(2)}(\mathbf{x}_k^{(2)}, \mathbf{u}_k)] \\ \mathbf{x}_{k+1}^{(3)} & = \mathbf{A}^{(3)}(\mathbf{x}_k^{(3)}) + \mathbf{K}^{(3)}(\mathbf{u}_k) \\ \mathbf{z}''_k & = \mathbf{C}_z^{(3)}(\mathbf{x}_k^{(3)}, \mathbf{u}_k) \\ \text{loss} & : \text{MSE}(\mathbf{z}''_k, \mathbf{C}_z^{(3)}(\mathbf{x}_k^{(3)}, \mathbf{u}_k)) \end{cases} \quad (\text{A.9})$$

We summarize BRAID’s three stages and their use-case in table A.1.

Table A.1: Summary of the three stages of BRAID and their use-case in learning various dynamics. Each stage has two RNNs: a predictor form and a generator form RNN, denoted without and with a *fw* subscript, respectively.

Stage	Models	Functionality
1	<i>RNN1, RNN1_{fw}</i>	Intrinsic behaviorally relevant neural dynamics
2	<i>RNN2, RNN2_{fw}</i>	Residual neural-specific intrinsic dynamics
3	<i>RNN3, RNN3_{fw}</i>	Residual behavior-specific dynamics not encoded in neural activity

A.1.4 MODEL ARCHITECTURE DETAILS AND HYPERPARAMETERS

Throughout the manuscript, to model nonlinearities within any of the transformations i.e., $\mathbf{A}(\cdot)$, $\mathbf{A}_{fw}(\cdot)$, $\mathbf{K}(\cdot)$, $\mathbf{K}_{fw}(\cdot)$, $\mathbf{C}_z(\cdot)$, and $\mathbf{C}_y(\cdot)$, we use a multi-layer perceptron (MLP), also known as a feedforward neural network, with a single hidden layer, 64 units in the hidden layer, and a *ReLU* nonlinearity as activation function. Otherwise, to keep a transformation linear, we replace the MLP with a linear mapping implementing a matrix multiplication, which is a special case of an MLP with no hidden layers and a linear activation function. For example, linear BRAID is a BRAID model with all mappings being linear whereas BRAID *Nonlinear* \mathbf{C}_z has a nonlinear MLP as the behavior decoder (both $\mathbf{C}_z^{(1)}$ and $\mathbf{C}_z^{(2)}$) while all of its other transformations are linear.

We use an Adam optimizer (Kingma & Ba, 2017) in all BRAID optimizations. We train models up to a maximum number of 2500 epochs to ensure convergence, while employing early stopping to avoid overfitting. We provide example learning curves in figure A.7. Details of the hyperparameters used for BRAID are provided in table A.2.

We also show that BRAID can locate the correct structure of nonlinearity within all possible combinations in our simulations (table 2). Here, we set each of the following four groups of transformations i.e., $\mathbf{A}(\cdot)/\mathbf{A}_{fw}(\cdot)$, $\mathbf{K}(\cdot)/\mathbf{K}_{fw}(\cdot)$, $\mathbf{C}_y(\cdot)$, $\mathbf{C}_z(\cdot)$, as linear or nonlinear, resulting in a total of 2^4 cases. To select one final configuration for the nonlinearity, we follow an automatic nonlinearity selection procedure for a given dataset. In this procedure, within the training data, we perform a 2-fold inner cross-validation in which we fit BRAID models with all 2^4 nonlinearity configurations and then pick the nonlinearity structure with the best cross-validated behavior decoding on the held-out section of the training data. Then we retrain a BRAID model with that selected structure on the

Table A.2: BRAID hyperparameters used in real data experiments and simulations

Hyperparameter	Value
Number of hidden layers in nonlinear maps	1
Number of hidden units in nonlinear maps	64
Nonlinear activation	ReLU
Learning rate	0.001
Batch size	32
Sequence length	128
Optimizer	Adam

entire training data to get our final model. Finally, we evaluate that final model on the unseen test data. We refer to this approach as automatic nonlinearity selection (table 2).

In simulations, state dimensions are set to be the same as that of the true model underlying the data. In real data analyses, to investigate the effect of n_x , we vary the state dimension in $n_x \in [1, 2, 4, 8, 16, 32, 64]$ and report the results (figure A.5). For BRAID and DPAD models, we always learn the first 16 dimensions via stage 1 i.e., $\mathbf{x}_k^{(1)} \in \mathbb{R}^{n_1}$ with $n_1 = \min(16, n_x)$, and if there is any more capacity left (i.e., if $n_x - n_1 = n_2$ is positive), it is dedicated to the irrelevant states $\mathbf{x}_k^{(2)} \in \mathbb{R}^{n_x - n_1}$ and is learned using stage 2. We pick 16 as the dimensionality of the behaviorally relevant states as in our experiments, because BRAID reached close to its peak behavior decoding at this dimension. We refer to models with state dimensions $n_x = 16$ and 64 as low and high-dimensional regimes, respectively.

A.1.5 INTRINSIC DYNAMICS AND EIGENVALUES

To evaluate how well the intrinsic behaviorally relevant neural dynamics are learned, in simulations (figures 2, A.1), we assess the eigenvalues of the generative transition \mathbf{A}_{fw} which characterize the intrinsic dynamics. Note that the ground truth and learned models in these simulations both have a linear intrinsic state transition \mathbf{A}_{fw} and the nonlinearity either lies in the transformation from latent states to behavior (figure 2) or transformation from external inputs to the latent space (figure A.1). We learn the forward recursion parameters \mathbf{A}_{fw} and \mathbf{K}_{fw} of equations 1 and 3 as we optimize multi-step-ahead predictions. We take the eigenvalues of the intrinsic transition \mathbf{A}_{fw} and compare them to that of the ground truth model. For the ground truth (λ_i) and identified ($\hat{\lambda}_i$) eigenvalues we first pair them as $\{\lambda_1, \lambda_2, \dots, \lambda_{n_1}\}$ and $\{\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_{n_1}\}$ such that the sum of squared distances of pairs is minimized. We then calculate the normalized eigenvalue error as:

$$\frac{\sqrt{\sum_{i=1}^{n_1} \|\lambda_i - \hat{\lambda}_i\|^2}}{\sqrt{\sum_{i=1}^{n_1} \|\lambda_i\|^2}}. \tag{A.10}$$

A.2 BASELINES

First, to assess the impact of the nonlinearities learned by BRAID, we compare it against two fully linear dynamical methods: IPSID (Vahidi et al., 2024) and linear BRAID. Second, to highlight the significance of modeling the effect of measured inputs on the neural-behavioral dynamics, we take an autonomous dynamical model, DPAD (Sani et al., 2024), as another baseline. Another important aspect of BRAID is supervision of the behaviorally relevant dynamics in presence of inputs in its first stage. To assess prioritization of the behaviorally relevant dynamics due to this supervision, we take an equivalent unsupervised baseline termed U-BRAID detailed below. We also compare BRAID against a multi-modal nonlinear method that allows accounting for external inputs termed mmPLRNN (Kramer et al., 2022). Finally, we compare BRAID to TNDM (Hurwitz et al., 2021), a second autonomous method for modeling neural-behavioral dynamics.

972 A.2.1 IPSID

973
974 IPSID, similar to BRAID, models the effect of measured external inputs on neural-behavioral dy-
975 namics but operates under a fully linear framework. It fits the parameters of a linear version of
976 equation A.1 via a projection-based analytical algorithm called subspace identification (Van Over-
977 schee & De Moor, 1996).

978 A.2.2 LINEAR BRAID

979
980 Linear BRAID serves as another linear baseline, retaining the same architecture and learning stages
981 as BRAID but with all transformations replaced by linear mappings. In essence, linear BRAID
982 and IPSID have a similar model that are learned differently. Linear BRAID uses the same numerical
983 optimization used in BRAID. BRAID reduces to linear BRAID by removing all hidden layers within
984 model transformations and setting all activation functions to linear. In simulations, we find that linear
985 BRAID and IPSID perform similarly as expected (figure 2).

986 A.2.3 DPAD

987
988 Dissociative Prioritized Analysis of Dynamics (DPAD) (Sani et al., 2024), learns a nonlinear model
989 that dissociates and prioritizes dynamics shared between neural activity and behavior, but it impor-
990 tantly does not account for the external inputs. Originally, DPAD also does not allow for multi-
991 step-ahead optimization and thus does not learn a generative form representation of the dynamics,
992 which is in contrast to the learning of A_{fw} in BRAID. We extend DPAD to add optimization of
993 multi-step-ahead predictions into the DPAD framework for a more fair comparison to BRAID in
994 terms of forecasting.

995 A.2.4 U-BRAID

996
997 U-BRAID is an unsupervised method, which only performs stage 2 of the BRAID learning proce-
998 dure. As such, U-BRAID learns all neural dynamics irrespective of their relevance to the behavior,
999 but still while considering inputs (equation A.5). U-BRAID does not utilize behavior information in
1000 learning dynamics, and the extracted latent states are later mapped to the behavior data via a down-
1001 stream decoder (equation A.6 but without $\mathbf{x}_k^{(1)}$). In fact, U-BRAID is special case of BRAID with
1002 $n_1 = 0$ and $n_x = n_2$.

1003 A.2.5 MMPLRNN

1004
1005 Multi-modal piecewise-linear RNN (mmPLRNN) is a method previously introduced for multi-
1006 modal dynamical modeling with piecewise-linear RNNs (Kramer et al., 2022). This method allows
1007 for modeling external inputs, although this aspect of it has not been investigated in any prior work.
1008 Nevertheless, we compare BRAID to an input-driven mmPLRNN to further assess its performance.
1009 mmPLRNN builds on a prior work, PLRNN (Durstewitz, 2017), by fusing information from two
1010 modalities (e.g., neural activity and behavior). By design, mmPLRNN utilizes both modalities (and
1011 input) during inference, which is in contrast to BRAID that only uses neural activity (and input)
1012 during inference. Although this provides the benefit of using more data for behavior decoding to
1013 mmPLRNN and confounds the comparison with BRAID, we still include the mmPLRNN results.
1014 We train mmPLRNN models with nonlinear readouts comparable to BRAID and compare their
1015 neural-behavioral forecasting. mmPLRNN is a generative model whose parameters are learned via
1016 variational inference. We used the recommended hyperparameters from the original work¹.

1017 A.2.6 TNDM

1018
1019 Targeted Neural Dynamical Modeling (TNDM) (Hurwitz et al., 2021) is a method based on sequen-
1020 tial autoencoders that learns two sets of dynamics: one contributing to both neural and behavioral
1021 data, and the other only contributing to neural data. TNDM uses a non-causal, bidirectional RNN as
1022 the encoder to infer the initial conditions for its relevant and irrelevant generator/decoder RNNs. The
1023 dynamical model is learned via variational inference with both neural and behavioral reconstructions
1024

1025 ¹We use the implementation provided in <https://github.com/DurstewitzLab/mmPLRNN>

1026 optimized simultaneously with a combined loss. Unlike BRAID, TNDM does not account for mod-
 1027 eling the effect of external inputs on neural-behavioral dynamics. Therefore, for comparisons, we
 1028 also implement an extension of TNDM to allow the inclusion of external inputs. In this version, we
 1029 provide the external inputs to TNDM model as input by concatenating them with the neural activity
 1030 as the input to the model. Importantly, we do not add reconstruction of inputs as part of the loss
 1031 to keep the loss the same as that of the original TNDM and keep the learned model focused on
 1032 neural-behavioral reconstruction.

1033 We compare BRAID to TNDM (with and without addition of sensory stimuli as external input) in
 1034 our real data experiments. Unlike BRAID that models the neural observations with a Gaussian distri-
 1035 bution, TNDM uses a Poisson observation model for the neural data. Therefore, in our comparisons
 1036 to TNDM, we analyze non-smoothed spike counts in 50ms bins (for both TNDM and BRAID). We
 1037 use the default hyperparameters from the original work² for TNDM.

1038 A.2.7 LFADS

1040 Latent Factor Analysis via Dynamical Systems (LFADS) (Pandarinath et al., 2018) is an unsuper-
 1041 vised method that combines nonlinear dynamical modeling with sequential autoencoders. Similar to
 1042 TNDM, LFADS uses a non-causal, bidirectional RNN as the encoder to infer the initial conditions
 1043 for a generator/decoder RNN. The dynamical model is learned via variational inference for unsu-
 1044 pervised neural reconstruction. LFADS has a version that uses additional RNNs called controller
 1045 networks to infer unmeasured inputs from the neural data and use those inferred inputs to drive its
 1046 generator RNN. We use this version of LFADS (O’Shea & Pandarinath, 2021)³, to serve as an addi-
 1047 tional benchmark in comparison to BRAID’s modeling of measured inputs. For a full visualization
 1048 of the LFADS model including the controller networks see Supplementary Fig. 12 in Pandarinath
 1049 et al., 2018. Briefly, without a controller network, all the information about the complete trial has
 1050 to be encoded into the initial state of the LFADS generator RNN, which then autonomously evolves
 1051 to extract states and factors over the course of the trial. In contrast, a controller network acts as a
 1052 regular non-autonomous RNN for LFADS and allows its generator RNN to take corrections from the
 1053 neural data throughout the trial, hence improving its capacity to accommodate non-smooth changes
 1054 in the middle of trials (Pandarinath et al., 2018). The controller network itself does not directly take
 1055 neural data as input, rather an additional bidirectional RNN called the controller-encoder operates
 1056 on the neural data of each trial, and the states of this controller-encoder are passed as input to the
 1057 controller network. As a result, the ‘inferred inputs’ in LFADS are also non-causally inferred. To
 1058 compare results with BRAID, we set the number of factors, and the latent states of the generator and
 1059 the controller-encoder networks to be the same as BRAID. We also pass the same smoothed neural
 1060 data to LFADS and use the Gaussian neural loss option accordingly. Other hyperparameters were
 1061 set as those in the original work (Pandarinath et al., 2018).

1062 A.2.8 CEBRA

1063 CEBRA (Schneider et al., 2023) is a recent method proposed for extracting latent embeddings from
 1064 neural data in a way that can be guided by behavior. CEBRA uses a 1-dimensional convolutional net-
 1065 work to extract embeddings from small windows of neural data and guides the extraction of latents
 1066 via a contrastive loss. The supervised version of CEBRA (i.e., CEBRA-Behavior) uses a contrastive
 1067 loss on behavior to learn embeddings that dissociate samples with different behavior data, and are
 1068 thus behaviorally relevant. Due to its use of a convolutional network to extract embeddings, CEBRA
 1069 does not learn explicit recursive dynamics and also extracts embedding from a finite window of neu-
 1070 ral data at a time. CEBRA also does not learn any models to decode behavior or neural data from
 1071 the extracted embeddings. Thus, as done in the original work (Schneider et al., 2023), we fit k-NN
 1072 regression⁴ models to decode neural-behavioral data from the extracted CEBRA embeddings. For
 1073 a fair comparison, we also provided the measured sensory inputs time-series of the task to CEBRA
 1074 by concatenating them with neural activity as the input. We use the default hyperparameters from
 1075 the original work⁵ for CEBRA.

1076 ²We use the implementation provided in <https://github.com/HennigLab/tndm>

1077 ³We use the implementation provided in <https://lfads.github.io/lfads-run-manager/>

1078 ⁴sklearn.neighbors.KNeighborsRegressor (Pedregosa et al., 2011)

1079 ⁵We use the implementation provided in [https://github.com/AdaptiveMotorControlLab/](https://github.com/AdaptiveMotorControlLab/CEBRA)
 CEBRA

1080 A.3 NON-HUMAN PRIMATE ELECTROPHYSIOLOGICAL RECORDINGS FROM

1081
1082 We analyzed a publicly available dataset (O’Doherty et al., 2017) in which a macaque (monkey
1083 I) performs a motor task. Spiking activity was recorded from primary motor cortex (M1), while
1084 the subject controlled a 2D cursor to reach targets that appeared on random locations on a grid
1085 within a virtual reality environment. Targets appeared back to back, without any time gaps. We
1086 took the subject’s 2D fingertip position and velocity as the behavior time-series \mathbf{z}_k , and the sensory
1087 input, taken as 2D location of the current target, as the input signal \mathbf{u}_k . We analyzed the first
1088 spike dimension available for each channel-resulting in 89 to 92 units from the first 7 available
1089 recording sessions and randomly selected half of these units to model as our neural activity. For
1090 neural modality, we use spike counts within 50 ms non-overlapping windows. Finally, we smoothed
1091 the spike counts by a Gaussian kernel with a 50 ms s.d. (except for in table A.5) and took that as the
1092 neural time-series \mathbf{y}_k . We report the mean and standard error of the mean (s.e.m.) computed across
1093 7 sessions and 5 cross-validated folds.

1094 As a second neural modality, we also model the raw local field potential (LFP) activity recorded
1095 from the same monkey during the same task. As the only preprocessing, we apply a 10 Hz anti-
1096 aliasing filter to be able to downsample the raw LFP to the sampling rate of behavior (i.e., 20 Hz).
1097 We refer to this modality as raw LFP data. Results for raw LFP data (figure A.3) were consistent
1098 with those obtained for spiking activity (figure 3).

1100 A.4 SIMULATION DETAILS

1101 We analyze three simulated datasets based on dynamical systems. In all the three, we generate 10
1102 different sets of random linear matrices for equation A.11, then generate the ground truth latent
1103 states \mathbf{x}_k , neural activity \mathbf{y}_k and behavior observations \mathbf{z}_k . In equation A.11, \mathbf{w}_k , \mathbf{v}_k , and ϵ_k are
1104 zero-mean white Gaussian noises accounting for unmeasured excitations, neural observation noise,
1105 and behavioral observation noise respectively. $f_{C_z}(\cdot)$ and $f_B(\cdot)$ are nonlinear functions, as described
1106 below for each simulation.
1107

$$1108 \begin{cases} \mathbf{x}_{k+1} &= \mathbf{A}_{fw}\mathbf{x}_k + f_B(\mathbf{B}\mathbf{u}_k) + \mathbf{w}_k \\ \mathbf{y}_k &= \mathbf{C}_y\mathbf{x}_k + \mathbf{D}_y\mathbf{u}_k + \mathbf{v}_k \\ \mathbf{z}_k &= f_{C_z}(\mathbf{C}_z\mathbf{x}_k) + \mathbf{D}_z\mathbf{u}_k + \epsilon_k \end{cases} \quad (\text{A.11})$$

1109
1110 To generate a temporally structured input in all cases, we simulate a separate random linear state
1111 space model according to equation A.12 and take its output as the external input \mathbf{u}_k to the main
1112 model of equation A.11.

$$1113 \begin{cases} \mathbf{x}_{k+1}^u &= \mathbf{A}_{fw}^u\mathbf{x}_k^u + \mathbf{w}_k^u \\ \mathbf{u}_k &= \mathbf{C}_u\mathbf{x}_k^u + \mathbf{v}_k^u \end{cases} \quad (\text{A.12})$$

1120 A.4.1 SIMULATION 1: SPIRAL BEHAVIOR MANIFOLD

1121 For the first simulation, we generate data with a spiral behavior manifold. To do so, we apply a
1122 pointwise nonlinear mapping $f_{C_z}(\nu) = \left[\frac{\bar{\nu}}{2} \cos(\bar{\nu}) \quad \frac{\bar{\nu}}{2} \sin(\bar{\nu}) \right]^T$ to the readout from the latent
1123 states (third line in equation A.11). The bar over the function input ν indicates a scaling factor
1124 that normalizes it before nonlinear function is applied. See figure 2a for a visualization of the
1125 nonlinearity. We take $f_B(\cdot)$ as identity, set dimensions to $n_y = n_z = n_u = n_x = n_1 = 2$, and take
1126 $\mathbf{D}_y = \mathbf{D}_z = 0$ for this simulation.
1127

1129 A.4.2 SIMULATION 2: TRIGONOMETRIC BEHAVIOR MANIFOLD

1130 For the second simulation with trigonometric behavior map, we apply another nonlinearity, point-
1131 wise sinusoidal nonlinear function, $f_{C_z}(\nu) = a \sin(\bar{\nu}) + b\bar{\nu}$, to the latent states to generate the
1132 behavior \mathbf{z}_k . See figure 2e for a visualization of the nonlinearity. In this simulation $f_B(\cdot)$ is taken as
1133 identity function and we set $n_y = n_z = n_u = n_x = n_1 = 1$.

1134 A.4.3 SIMULATION 3: TRIGONOMETRIC INPUT-ENCODER

1135 For the third simulation, as we iterate over the state equation (first line in equation A.11), we apply
 1136 a pointwise sinusoidal nonlinear function, $f_B(\nu) = a \sin(\bar{\nu}) + b\bar{\nu}$, to the input. See figure A.1a for a
 1137 visualization of the nonlinearity. Here $f_{C_z}(\cdot)$ is taken to be an identity function. In this simulation,
 1138 we set $n_y = n_z = n_u = n_x = n_1 = 1$.
 1139

1140 We also perform two additional simulations (figure A.2) that are similar in nonlinearity structure to
 1141 the second and third simulations explained above, but incorporate an additional 1-dimensional latent
 1142 state $\mathbf{x}_k^{(3)}$, as in figure 1a, representing input-driven behavior-specific dynamics not encoded in the
 1143 neural activity, as follows:

$$\begin{cases} \begin{bmatrix} \mathbf{x}_{k+1}^{(1)} \\ \mathbf{x}_{k+1}^{(3)} \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_{f_w}^{(1)} \mathbf{x}_k^{(1)} \\ \mathbf{A}_{f_w}^{(3)} \mathbf{x}_k^{(3)} \end{bmatrix} + \begin{bmatrix} f_B(\mathbf{B}^{(1)} \mathbf{u}_k) \\ \mathbf{B}^{(3)} \mathbf{u}_k \end{bmatrix} + \mathbf{w}_k \\ \mathbf{y}_k &= \mathbf{C}_y^{(1)} \mathbf{x}_k^{(1)} + \mathbf{D}_y \mathbf{u}_k + \mathbf{v}_k \\ \mathbf{z}_k &= f_{C_z}(\mathbf{C}_z^{(1)} \mathbf{x}_k^{(1)}) + \mathbf{C}_z^{(3)} \mathbf{x}_k^{(3)} + \mathbf{D}_z \mathbf{u}_k + \epsilon_k \end{cases} \quad (\text{A.13})$$

1151 A.5 SUPPLEMENTARY RESULTS

1153 A.5.1 BRAID CAN EXCLUDE NON-ENCODED BEHAVIOR-SPECIFIC DYNAMICS

1154 Here, we demonstrate that the optional preprocessing step in BRAID (detailed in section 3.3) can
 1155 dissociate behavior-specific dynamics (i.e., those that are not encoded in the neural activity) during
 1156 learning and make sure they are not conflated with intrinsic neural dynamics and are not mixed into
 1157 the neural states ($\mathbf{x}_k^{(1)}$ and $\mathbf{x}_k^{(2)}$). We conducted two additional simulations similar in structure to the
 1158 second and third simulations explained in section 4.1.2. However, here, for all simulated models,
 1159 we added input-driven dynamics that influenced behavior but were not encoded in neural activity
 1160 (denoted as $\mathbf{x}_k^{(3)}$ in figure 1a and equation A.13). The preprocessing step is intentionally expected to
 1161 yield latent states that are potentially less predictive of behavior, but are encoded in neural activity.
 1162 When desired, BRAID provides the option to further learn behavior-specific dynamics post-hoc
 1163 with a separate latent state ($\mathbf{x}_k^{(3)}$). The preprocessing and post-hoc learning steps allow BRAID to
 1164 avoid conflation of non-encoded behavior dynamics with others, while also being able to learn these
 1165 dynamics and thus not incurring any overall reduction in behavior decoding.
 1166

1167 We fitted BRAID models with the preprocessing, and both with and without post-hoc learning of
 1168 behavior-specific dynamics. With the preprocessing, BRAID reached the neural prediction perfor-
 1169 mance of the ground truth model indicating correct removal of behavior-specific dynamics (figure
 1170 A.2). Moreover, the optional learning of behavior-specific dynamics led to reaching the behavior
 1171 decoding performance of the ground truth model (figure A.2), suggesting that one could option-
 1172 ally learn these dynamics as well within BRAID to gain interpretability (by learning a disentangled
 1173 model) without compromising decoding performance.

1174 A.5.2 ADDITIONAL SUPPLEMENTARY TABLES

1175 In this section we include additional supplementary tables that further support the results from the
 1176 main text. The caption for each supplementary table includes all the details, but here we provide a
 1177 list of these tables:
 1178

- 1179 • Table A.3: Ablation analysis showing the importance of the RNN_{f_w} generative model in BRAID
 1180 for learning intrinsic dynamics.
- 1181 • Table A.4: BRAID results for different nonlinearity configurations in real NHP data.
- 1182 • Table A.5: Comparison with TNM in real NHP data.
- 1183 • Table A.6: Comparison with LFADS and CEBRA in real NHP data.
- 1184 • Table A.7: BRAID results for modeling different number of neurons in real NHP data.
- 1185 • Table A.8: Same as table 3 shown in terms of the R^2 metric.

Table A.3: **A separate generative model is essential to learn the intrinsic dynamics accurately with BRAID.**

Row 1: BRAID, when learning a separate generative model (RNN_{fw}), more accurately learns the intrinsic dynamics as quantified by the error in identifying eigenvalues of the ground truth intrinsic dynamics in the simulated dataset with spiral manifold in figure 2a.

Row 2: The error when ablating the forward RNN from BRAID.

Rows 3-4: DPAD, even when optimized with m -step-ahead prediction loss and/or with input (\mathbf{u}_k), does not learn the intrinsic dynamics accurately.

Method	\log_{10} normalized eigenvalue error
BRAID	-1.3963 \pm 0.2551
BRAID without RNN_{fw}	0.0635 \pm 0.2212
DPAD + m -step loss	0.0357 \pm 0.1587
DPAD + input (\mathbf{u}_k) + m -step loss	-0.6512 \pm 0.0354

Table A.4: Comparison of BRAID model’s nonlinearity configurations in the NHP dataset ($n_x = n_1 = 16$, 4-step-ahead).

Model nonlinearity	Behavior forecasting CC	Neural forecasting CC
Linear	0.7453 \pm 0.0066	0.1767 \pm 0.0054
Recursion (\mathbf{A} , \mathbf{A}_{fw})	0.7121 \pm 0.0059	0.2719 \pm 0.0061
Encoder (\mathbf{K} , \mathbf{K}_{fw})	0.7181 \pm 0.0078	0.1646 \pm 0.0049
Decoder (\mathbf{C}_z , \mathbf{C}_y)	0.8042 \pm 0.0085	0.3274 \pm 0.0078

Table A.5: Comparison to TNDM, both when sensory input is additionally provided to the TNDM model and when it is not (see appendix A.2.6). All models are learned in low-dimensional regime i.e., BRAID with $n_x = n_1 = 16$, and TNDM with 16 relevant factors only. We used non-smoothed spike counts as the neural signals in this analysis. BRAID performances are for causal 1-step-ahead prediction, whereas the TNDM performances are non-causal smoothing performances, which are the only option for TNDM since it is a sequential autoencoder.

Method	Behavior decoding CC	Neural prediction CC
TNDM	0.3752 \pm 0.0170	0.3021 \pm 0.0051
TNDM with sensory input	0.6219 \pm 0.0103	0.3075 \pm 0.0050
BRAID (ours)	0.7841 \pm 0.0079	0.2935 \pm 0.0053

Table A.6: Comparison to LFADS (with controller), and CEBRA (CEBRA-behavior). LFADS inferred the external input to the dynamical system (appendix A.2.7). For CEBRA, w.s.i. (with sensory input) indicates that sensory input is additionally provided to the model to obtain the embeddings (appendix A.2.6). All models are learned with both low-dimensional ($n_x = 16$) and high-dimensional ($n_x = 64$) latent states ($n_1 = 16$ for BRAID). BRAID performances are for causal 1-step-ahead prediction, whereas LFADS and CEBRA performances are non-causal smoothing, and 0-step-ahead reconstruction respectively.

Method	Behavior decoding CC		Neural prediction CC	
	$n_x = 16$	$n_x = 64$	$n_x = 16$	$n_x = 64$
LFADS	0.4714 \pm 0.0192	0.5891 \pm 0.0135	0.5615 \pm 0.0086	0.5920 \pm 0.0072
CEBRA w.s.i.	0.7544 \pm 0.0073	0.7514 \pm 0.0072	0.5070 \pm 0.0053	0.5693 \pm 0.0041
BRAID (ours)	0.8109 \pm 0.0074	0.8085 \pm 0.0076	0.5571 \pm 0.0051	0.8401 \pm 0.0061

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

Table A.7: Effect of the number of neurons on BRAID’s performance. Results of the BRAID modeling when different numbers of neurons included as the neural signal. We included neurons from the same channel sets across different sessions and the ranges indicate the number of neurons available from those channels across different sessions. In all 3 rows, neural predictions are evaluated on the same common set neurons (i.e., the smallest set shown in row 1) to make the neural forecasting results comparable across rows. Behavior forecasting improved with more neurons and neural forecasting remained largely stable. These results suggest that BRAID can aggregate behaviorally relevant information across larger populations of neurons, while still being able to model this higher-dimensional population activity well.

Scale	Behavior forecasting CC		Neural forecasting CC	
	$n_x = 16$	$n_x = 64$	$n_x = 16$	$n_x = 64$
20-21 neurons	0.7727 ± 0.0091	0.7709 ± 0.0089	0.3206 ± 0.0081	0.4115 ± 0.0094
41-43 neurons	0.8042 ± 0.0085	0.7970 ± 0.0086	0.3220 ± 0.0088	0.4202 ± 0.0091
89-92 neurons	0.8337 ± 0.0061	0.8302 ± 0.0072	0.3329 ± 0.0082	0.4195 ± 0.0088

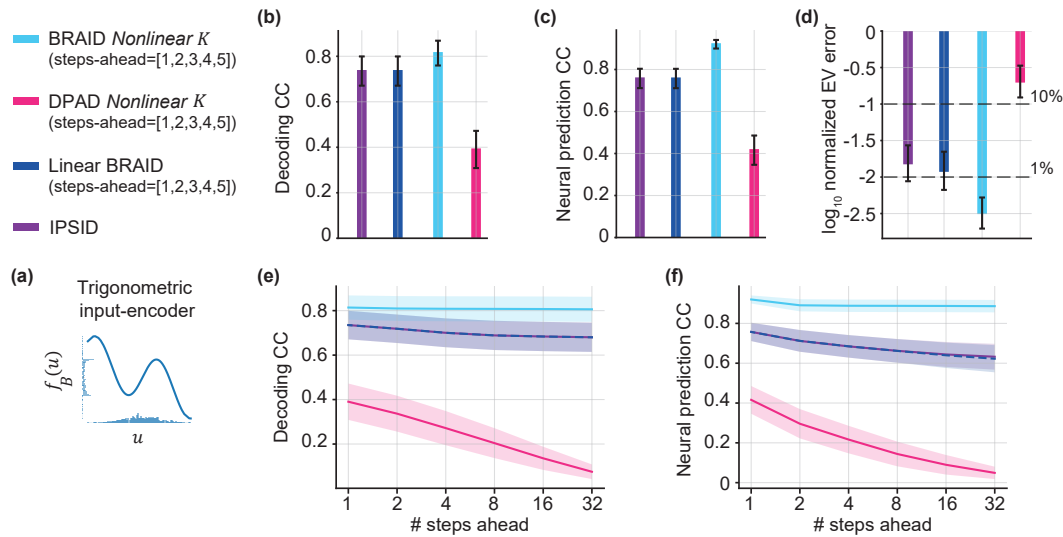
Table A.8: R^2 results for the same analyses provided in table 3. Forecasting performance (4-step-ahead) compared to baselines in NHP dataset for models with low-dimensional ($n_x = 16$) and high-dimensional ($n_x = 64$) latent states. $n_1 = 16$ for BRAID, linear BRAID, and DPAD. 1 of the 7 sessions were excluded from this table due to a very negative outlier in the mmPLRNN results.

Method	Behavior forecasting R^2		Neural forecasting R^2	
	$n_x = 16$	$n_x = 64$	$n_x = 16$	$n_x = 64$
linear BRAID	0.5821 ± 0.0056	0.5763 ± 0.0084	0.0239 ± 0.0036	0.1519 ± 0.0067
DPAD	0.4561 ± 0.0141	0.5497 ± 0.0101	0.0321 ± 0.0042	0.1344 ± 0.0069
U-BRAID	0.6047 ± 0.0097	0.6684 ± 0.0088	0.1768 ± 0.0068	0.1860 ± 0.0064
mmPLRNN	0.4737 ± 0.0118	0.6127 ± 0.0369	0.0734 ± 0.0079	0.0563 ± 0.0670
BRAID (ours)	0.6680 ± 0.0118	0.6578 ± 0.0129	0.1083 ± 0.0060	0.1792 ± 0.0064

1296 A.5.3 ADDITIONAL SUPPLEMENTARY FIGURES
 1297

1298 In this section we include supplementary figures that further support the results from the main text.
 1299 The caption for each supplementary figure include all the details, but here we provide a list of these
 1300 supplementary figures:

- 1301 • Figure A.1: Simulation results for simulation 3 with a trigonometric nonlinear input-encoder.
- 1302 • Figure A.2: Validation of the optional third stage of learning in BRAID for learning behavior-
- 1303 specific input driven dynamics.
- 1304 • Figure A.5: Results in real NHP data for different latent state dimensions.
- 1305 • Figure A.4: Average low-dimensional latent state trajectory extracted from real NHP data.
- 1306 • Figure A.6: Example BRAID decoded time series in real NHP data.
- 1307 • Figure A.3: Results of modeling the raw LFP modality in the real NHP data.
- 1308 • Figure A.7: Example plot showing loss versus epochs for BRAID.
- 1309
- 1310
- 1311



1330 **Figure A.1: BRAID results, optimized for forecasting, in simulation with trigonometric input-**
 1331 **encoder. (a)** Visualization of example nonlinearity in the simulation. **(b-c)** 1-step-ahead behavior
 1332 decoding and neural prediction for nonlinear BRAID, nonlinear DPAD, linear BRAID, and IPSID.
 1333 **(d)** Error in identifying intrinsic dynamics of the true model, quantified by the eigenvalues of the
 1334 state transition matrix \mathbf{A}_{f_w} . **(e-f)** Behavior and neural forecasting accuracy for 1 to 32 steps ahead,
 1335 enabled by learning the intrinsic dynamics (\mathbf{A}_{f_w}), with predictions optimized for [1, 2, 3, 4, 5]-steps-
 1336 ahead (section 3.1).
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

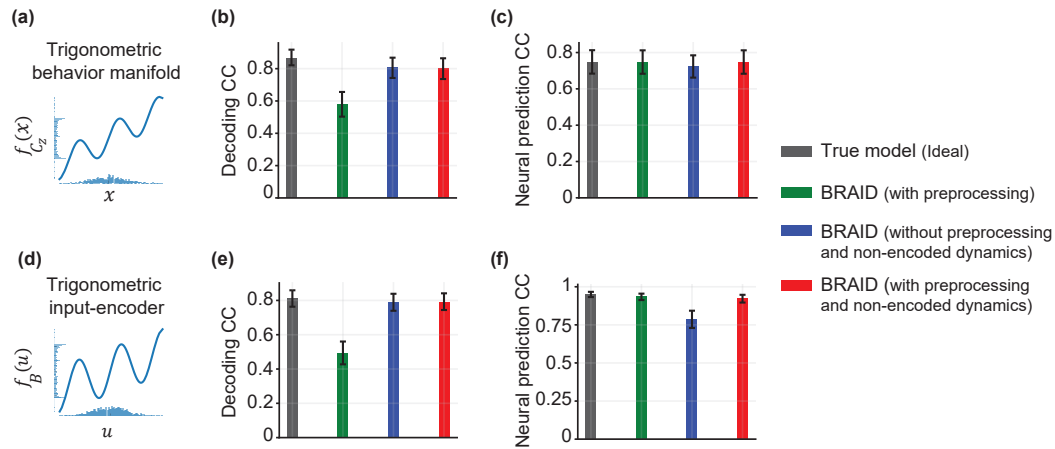


Figure A.2: **Behavior preprocessing successfully excludes non-encoded, behavior-specific dynamics in simulations.** 1-step-ahead behavior decoding and neural predictions for simulations with (a-c) trigonometric behavior decoder, and (d-f) input-encoder. Note that the true model includes behavior-specific dynamics, so here we expect that decoding of BRAID with preprocessing but without the post-hoc learning of behavior-specific dynamics (shown as green), to be worse than that of true model. Once the post-hoc learning step is also performed (shown as red), BRAID reaches ideal performance, but importantly does so while these behavior-specific dynamics are dissociated into a separate latent state $C_z^{(3)}$. In contrast, without the preprocessing step (shown as blue), BRAID reaches ideal decoding performance, but does so without having dissociated behavior specific dynamics to not be included in $C_z^{(1)}$. See section A.1.3 for details.

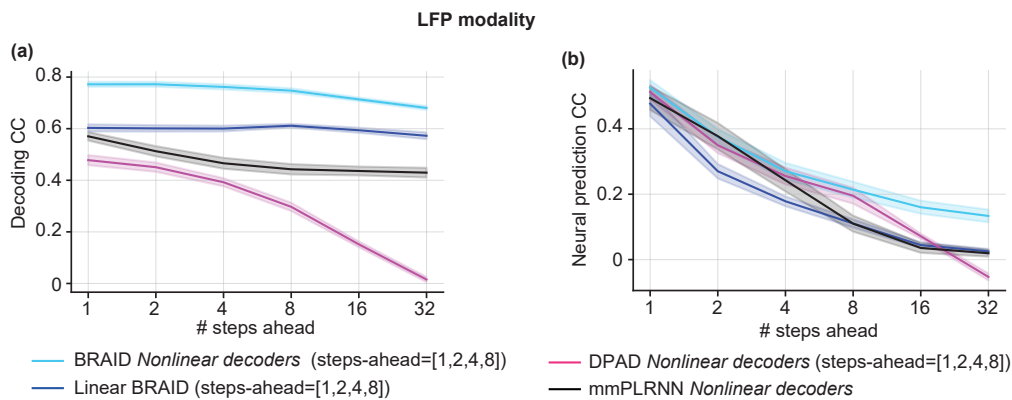
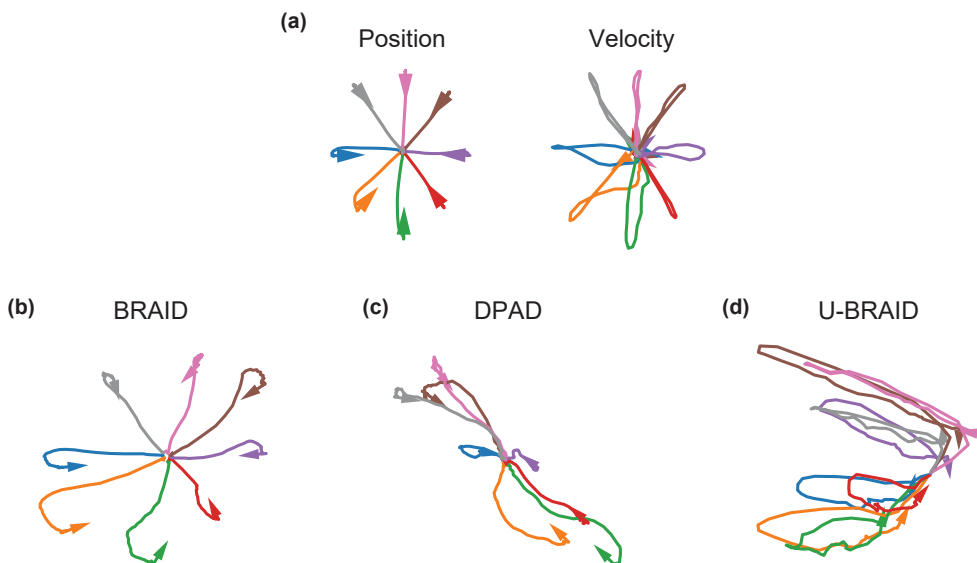


Figure A.3: **Analysis of Local Field Potential (LFP) neural modality. BRAID outperforms baselines in neural-behavioral forecasting.** We used LFP neural data as a second different modality and performed analysis similar to the one in figure 3 for smoothed spike counts. (a) Behavior and, (b) neural activity forecasting correlation coefficient (CC) for BRAID, linear BRAID, DPAD, and mmPLRNN for $n_x=16$. Shaded areas show the s.e.m., across the 3 recording sessions and 5 cross-validation folds. NHP results in all other figures and tables are for spiking data.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425



1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448

Figure A.4: **Latent states trajectories revealed on non-human primate reaching dataset.** We divide reach trials to 8 conditions based on reach direction (shown by colors) and find the condition-averaged (4-step-ahead) latent states trajectories for 2 dimensional models. (a) Condition-averaged behavior i.e., movement position and velocity. (b-d) Condition-averaged latent state trajectories for (b) BRAID, (c) DPAD and (d) U-BRAID. BRAID learns the most well-separated posterior (latent) trajectories for the 8 conditions in the task. U-BRAID’s trajectories are the least separated, showing that BRAID is more successful in extracting the behaviorally relevant intrinsic dynamics that are more congruent with behavior. Also, DPAD’s trajectories are not as well-separated as BRAID’s, although they are more separated than the unsupervised version (U-BRAID) as expected.

1449
1450
1451
1452
1453
1454
1455
1456
1457

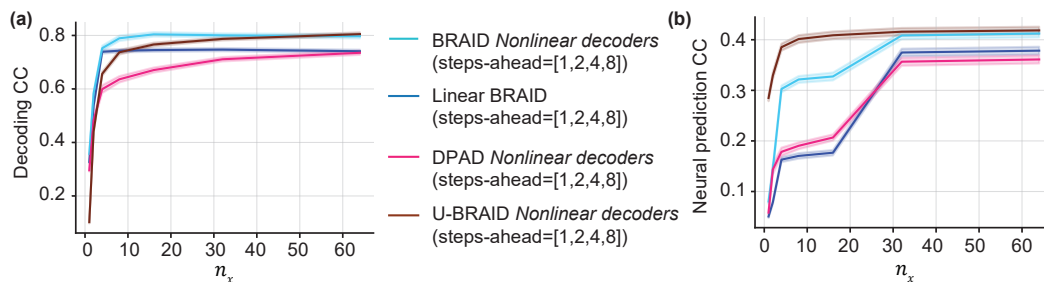
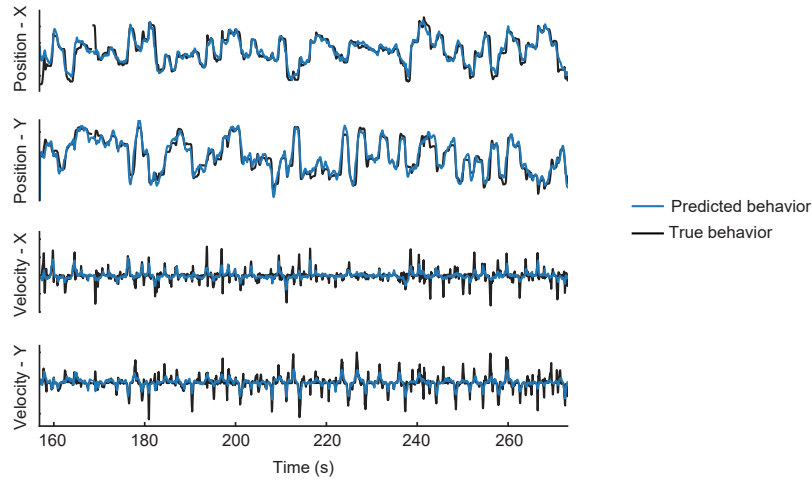


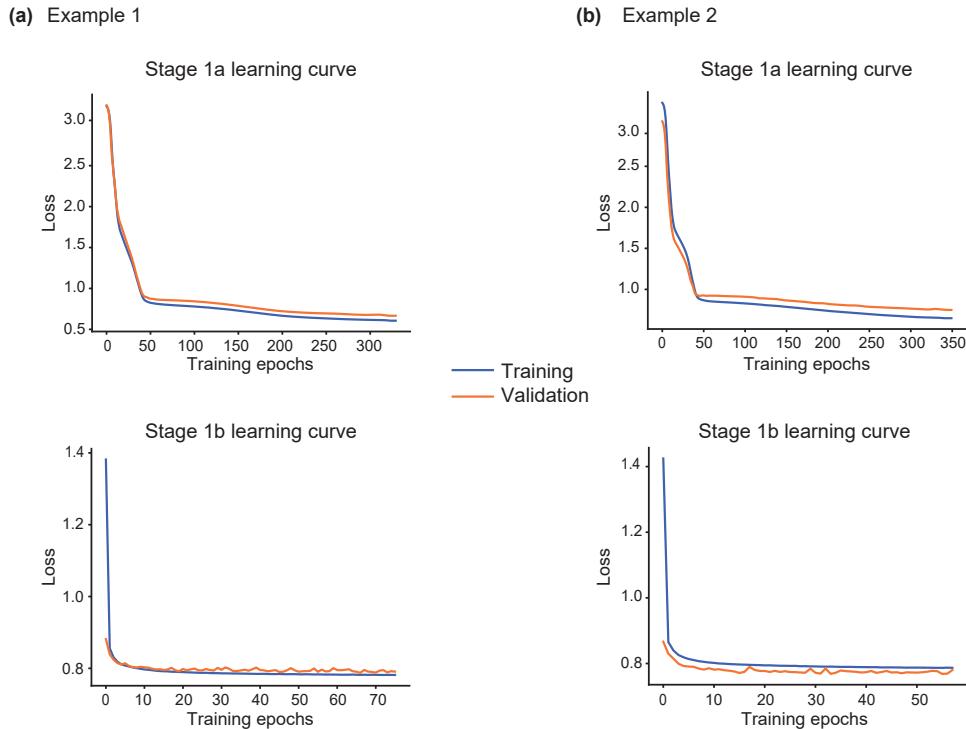
Figure A.5: **Forecasting (4-step-ahead predictions) correlation coefficient across latent dimensions.** BRAID predicts both (a) behavior and (b) neural activity more accurately than DPAD due to modeling input, and than linear BRAID due to modeling nonlinearity. As state dimension increases beyond 16 (dedicated to behaviorally relevant dynamics, i.e., $n_1 = 16$), BRAID uses stage 2 to learn neural specific dynamics, thus reaching the unsupervised baseline, i.e., U-BRAID, in neural prediction. For BRAID and DPAD, the first 16 state dimensions are dedicated to the behaviorally relevant neural dynamics (i.e., $n_1 = 16$) while any remaining dimensions ($n_x > 16$) are dedicated to the residual non-shared neural dynamics. 4-steps-ahead in this dataset corresponds to 200ms ahead.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475



1476
1477 **Figure A.6: Predicted behavior visualization.** True behavior versus BRAID’s (4-step-ahead) pre-
1478 dicted behavior for a representative session corresponding to the results in table 3 ($n_x = 64$).
1479

1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506



1507 **Figure A.7: Learning curve examples for BRAID model.** Loss values as a function of number
1508 of epochs for training and validation datasets for two example runs of BRAID. Top row: learning
1509 curve examples for the $RNNI$, $RNNI_{fw}$, i.e., stage 1a, in section 3.2. Bottom row: learning curves
1510 examples for stage 1b in section 3.2. These examples are taken from the same models whose per-
1511 formances are reported in figure 3 and table 3.