

# Extraction and Recovery of Spatio-Temporal Structure in Latent Dynamics Alignment with Diffusion Model (ERDiff)



NeurIPS'23 Spotlight

Authors: Yule Wang, Zijing Wu,  
Chengrui Li, Anqi Wu

Presenter: Yule Wang, ML@GT

Advisor: Prof. Anqi Wu

10/19/2023

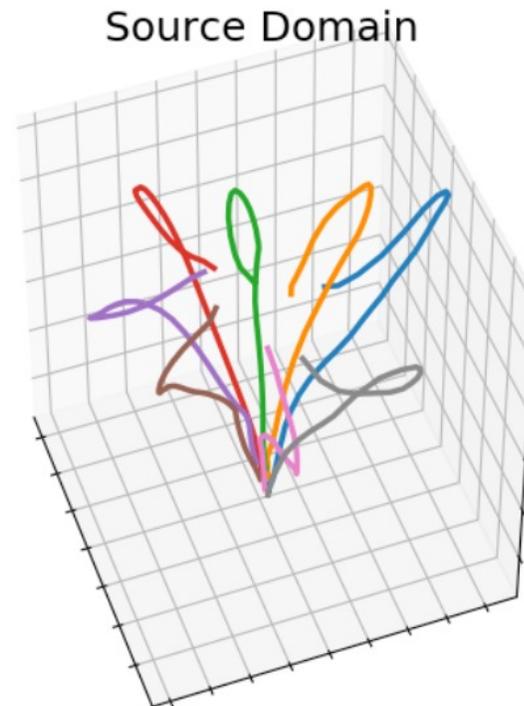
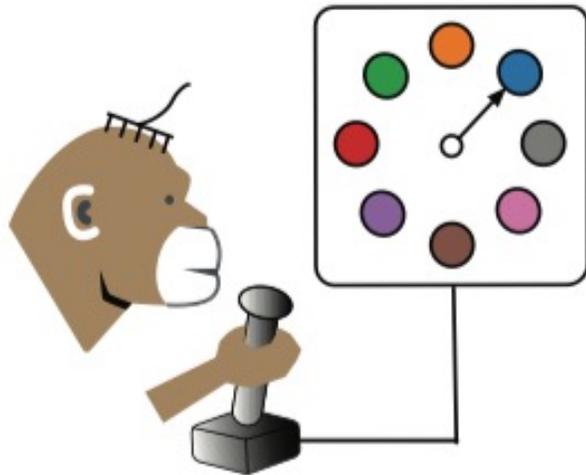


# Problem Formulation: Neural Distribution Alignment

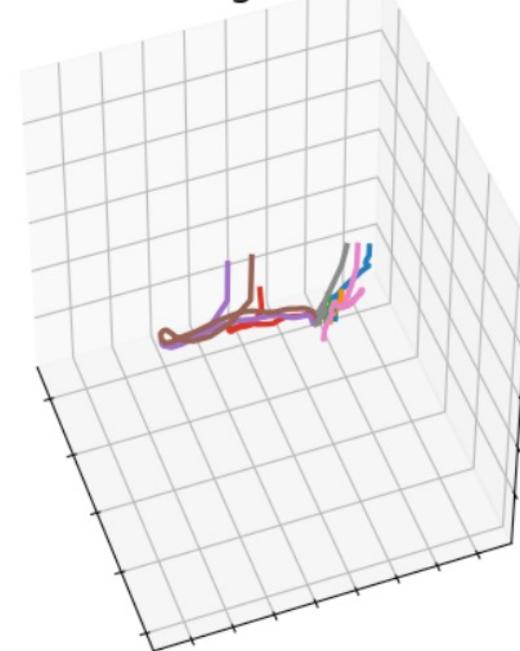
# Neural Distribution Alignment

## Problem to Solve:

- Perform alignment against the drastic distribution shift between raw neural population activities recorded across-day and inter-subject (individual).
- Task Setting Example
- Neural Latent Dynamics of Single Trials



Unaligned Target Domain  
Decoding  $R^2$ : -34.12



# Neural Distribution Alignment

## Problem Formulation:

### Distribution Alignment in low-dimensional latent space:

- Based on the Manifold hypothesis [1]: a relatively small number of latent dynamical factors contain a large portion of neural activities variability
- Perform alignment in the neural latent space

# Neural Distribution Alignment

## Problem Formulation:

### Source-Domain

- We denote the source-domain single-trial neural population activities as  $\mathbf{X}^{(s)} = [\mathbf{x}_1^{(s)}, \dots, \mathbf{x}_l^{(s)}]^\top \in \mathbb{R}^{l \times n}$ , where  $l$  is the trial length and  $n$  is the number of neurons.  $\mathbf{Z}^{(s)} = [\mathbf{z}_1^{(s)}, \dots, \mathbf{z}_l^{(s)}]^\top \in \mathbb{R}^{l \times d}$  are their low-dimensional latent dynamics inferred by a latent variable model (LVM), where  $d$  is the latent dimension size.
- $\phi_s, \psi_s = \operatorname{argmax}_{\phi, \psi} \left[ \mathbb{E}_{q(\mathbf{Z}^{(s)} | \mathbf{X}^{(s)}; \phi)} [\log p(\mathbf{X}^{(s)} | \mathbf{Z}^{(s)}; \psi)] - \mathbb{D}_{\text{KL}} \left( q(\mathbf{Z}^{(s)} | \mathbf{X}^{(s)}; \phi) \parallel \bar{q}(\mathbf{Z}^{(s)}) \right) \right]$
- We denote the probabilistic encoder and decoder of the LVM trained on source-domain as  $q(\mathbf{Z}^{(s)} | \mathbf{X}^{(s)}; \phi_s)$  and  $p(\mathbf{X}^{(s)} | \mathbf{Z}^{(s)}, \psi_s)$ , respectively.

# Neural Distribution Alignment

## Problem Formulation:

### Target-Domain

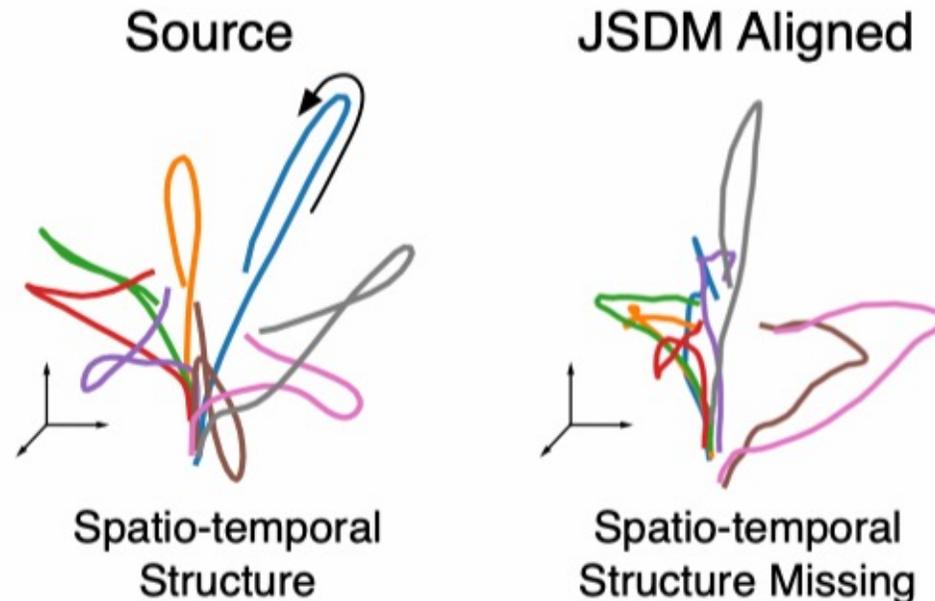
- Given the target-domain single-trial neural population activities as  $\mathbf{X}^{(t)} = [\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_l^{(t)}]^\top \in \mathbb{R}^{l \times n}$ , we perform distribution alignment by probing the probabilistic encoder  $q(\mathbf{Z}^{(t)} | \mathbf{X}^{(t)}; \boldsymbol{\phi})$
- The alignment is conducted by minimizing certain probability divergence  $\mathbb{D}(\cdot|\cdot)$  between the two posterior distributions:

$$\min_{\boldsymbol{\phi}_t} \mathbb{D} \left( q(\mathbf{Z}^{(s)} | \mathbf{X}^{(s)}; \boldsymbol{\phi}_s) \parallel q(\mathbf{Z}^{(t)} | \mathbf{X}^{(t)}; \boldsymbol{\phi}_t) \right)$$

# Neural Distribution Alignment

## Challenges:

- The underlying spatio-temporal structures of neural latent dynamics are both non-linear and complex.
- Prior works ignoring such crucial spatio-temporal structure information, resulting in comparably inferior alignment performance.



# Methodology of ERDiff

# Neural Distribution Alignment

## Key Insight:

- Extraction and Recovery of Spatio–Temporal Structure in Latent Dynamics Alignment with **Diffusion Model (ERDiff)**
- How to precisely extract the spatio–temporal structures of the source domain neural latent dynamics?

Strong Density Estimation Capability

- How to precisely recover such extracted spatio–temporal structures to the aligned target domain neural latent dynamics?

Proper Alignment Algorithm

# Neural Distribution Alignment

## Key Insight:

- How to precisely extract the spatio-temporal structures of the source domain neural latent dynamics?

Strong Density Estimation Capability

(Score-based) Diffusion Model

- How to precisely recover such extracted spatio-temporal structures to the aligned target domain neural latent dynamics?

Proper Alignment Algorithm

# Methodology of ERDiff

## Source-domain Diffusion Model (DM) Training:

- The LVM alone focuses on the point-to-point underlying mapping between  $\mathbf{X}^{(s)}$  and  $\mathbf{Z}^{(s)}$ . The overall distribution of the source domain latent dynamics  $p_s(\mathbf{Z})$ , abbreviation for  $q(\mathbf{Z}^{(s)} | \mathbf{X}^{(s)}; \phi_s)$ , is still inaccessible by building a LVM alone
- We propose to learn  $p_s(\mathbf{Z})$  via a Diffusion Model (DM) by taking the entire trials of latent dynamics on source-domain  $\mathbf{Z}^{(s)} \sim q(\cdot | \mathbf{X}^{(s)}; \phi_s)$  as input to the DM for training
- Specifically, the DM fits  $p_s(\mathbf{Z})$  through training of denoiser  $\epsilon(\mathbf{Z}, t; \theta_s): (\mathbb{R}^{l \times d} \times \mathbb{R}) \rightarrow \mathbb{R}^{l \times d}$  through the denoising score matching (DSM) loss:

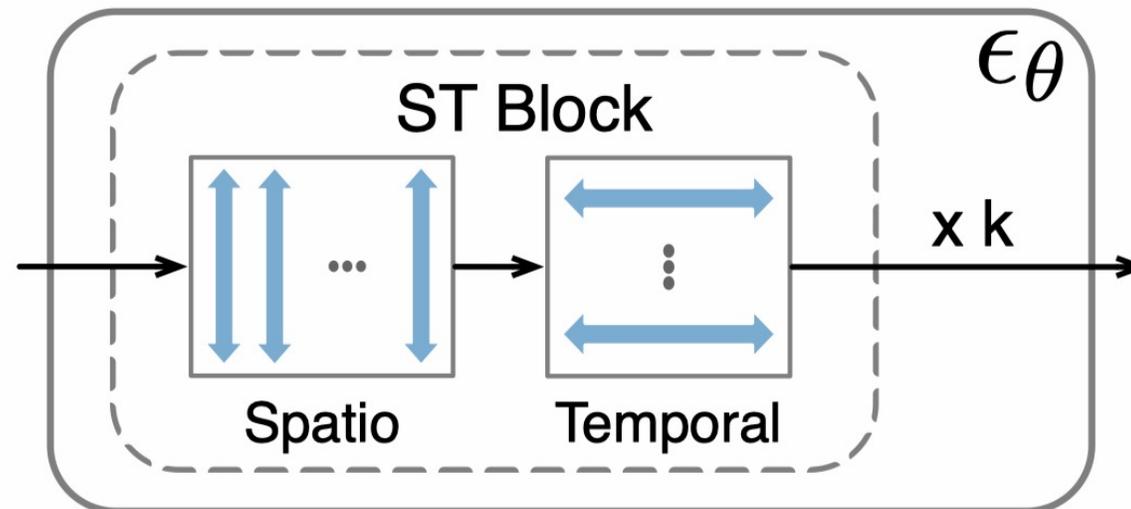
- $$\theta_s = \operatorname{argmin}_{\theta} \mathbb{E}_{t \sim \mathcal{U}[0, T]} \mathbb{E}_{\mathbf{Z}_0^{(s)} \sim q(\cdot | \mathbf{X}^{(s)}; \phi_s), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{l \times d})} \left[ w(t)^2 \left\| \epsilon - \epsilon(\mathbf{Z}_t^{(s)}, t; \theta) \right\|_2^2 \right]$$

- where  $\mathbf{Z}_t = \mu_t \mathbf{Z}_0 + \mathbf{K}_t \epsilon$ , and  $\epsilon(\mathbf{Z}_t, t; \theta) = -\mathbf{K}_t^{-T} \mathbf{s}(\mathbf{Z}_t, t; \theta)$ ,  $\mathbf{K}_t \mathbf{K}_t^T = \Sigma_t$

# Methodology of ERDiff

## Extract the spatio-temporal structures of the source domain:

- We use Spatio-Temporal Transformer Blocks (STBlocks)
- Each STBlock is composed of Spatio-Transformer layers followed by Temporal-Transformer layers
- Spatio-Transformer layer takes latent states of each time bin as inputs to extract spatial structure
- Temporal-Transformer layer takes the entire latent trajectory of each latent space dimension as inputs to extract temporal structure



# Neural Distribution Alignment

## Key Insight:

- How to precisely extract the spatio-temporal structures of the source domain neural latent dynamics?

Strong Density Estimation Capability

- How to precisely recover such extracted spatio-temporal structures to the aligned target domain neural latent dynamics?

Proper Alignment Algorithm

**Maximum Likelihood Alignment** with Diffusion Model

# Methodology of ERDiff

## Maximum Likelihood Alignment with Diffusion Model:

- Given the target-domain neural activities  $\mathbf{X}^{(t)}$ , we propose to perform distribution alignment via maximum likelihood estimation (MLE):

$$\operatorname{argmax}_{\phi} \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X}^{(t)})} [\log p_s(q(\mathbf{Z} | \mathbf{X}; \phi))] = \operatorname{argmax}_{\phi} \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}^{(t)}; \phi)} [\log p_s(\mathbf{Z})],$$

- We use the marginal distribution  $p_0(\mathbf{Z}; \theta_s)$  learnt by the DM to approximate  $p_s(\mathbf{Z})$ , the maximum likelihood estimation can thus be written as:

$$\operatorname{argmax}_{\phi} \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}^{(t)}; \phi)} [\log p_0(\mathbf{Z}; \theta_s)]$$

- We note that the alignment is conducted by probing the parameter set  $\phi$  of the probabilistic encoder while keeping the DM  $p_0(\mathbf{Z}; \theta_s)$  fixed

# Methodology of ERDiff

## Maximum Likelihood Alignment with Diffusion Model:

- Consequently, we could obtain an upper bound of the maximum likelihood loss function, as follows:

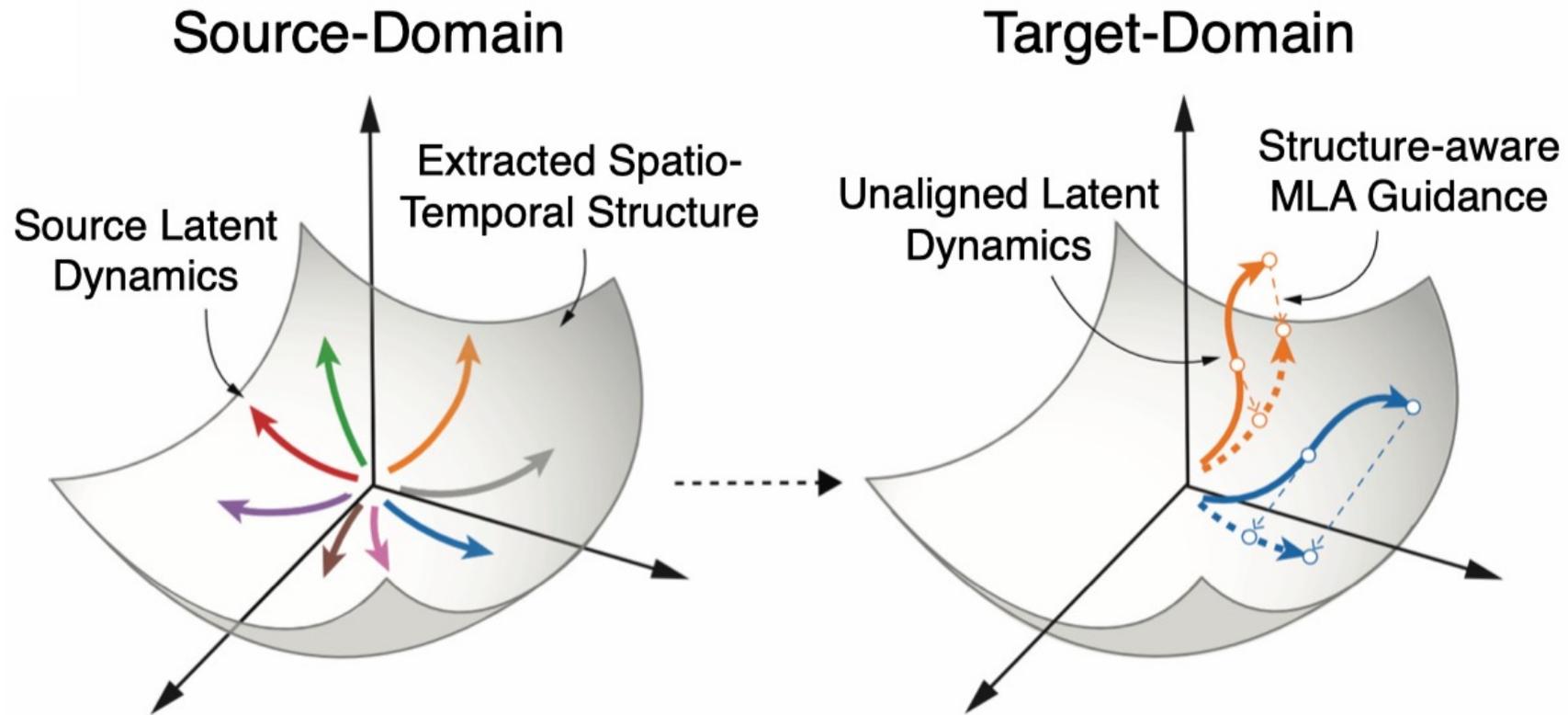
$$\begin{aligned} -\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X}^{(t)}; \phi)} [\log p_0(\mathbf{Z}; \theta_s)] &\leq \underbrace{\mathbb{D}_{\text{KL}}(p_T(\mathbf{Z}; \theta_s) \parallel \pi(\mathbf{Z}))}_{\text{Constant Term}} \\ &+ \mathbb{E}_{t \sim \mathcal{U}[0, T]} \mathbb{E}_{\mathbf{Z}_0 \sim q(\mathbf{Z}|\mathbf{X}^{(t)}; \phi), \epsilon \sim \mathcal{N}(0, I_{l \times d})} \left[ \underbrace{w(t)^2 \|\epsilon - \epsilon(\mathbf{Z}_t, t; \theta_s)\|_2^2}_{\text{Denoising Score Matching}} - \underbrace{2\nabla_{\mathbf{Z}} \cdot \mathbf{f}(\mathbf{Z}_t, t)}_{\text{Divergence}} \right] \end{aligned}$$

- The divergence objective can be approximated using the Hutchinson–Skilling trace estimator [1], making the whole optimization objective computationally tractable.
- We note that the recovery of spatio-temporal structure is primarily conducted by the denoising score matching part.

[1] J. Skilling, “The eigenvalues of mega-dimensional matrices,”

# ERDiff Spatio-Temporal Structure Recovery

## Recovery



# Experiment

# Experiment

## Synthetic Datasets – Results

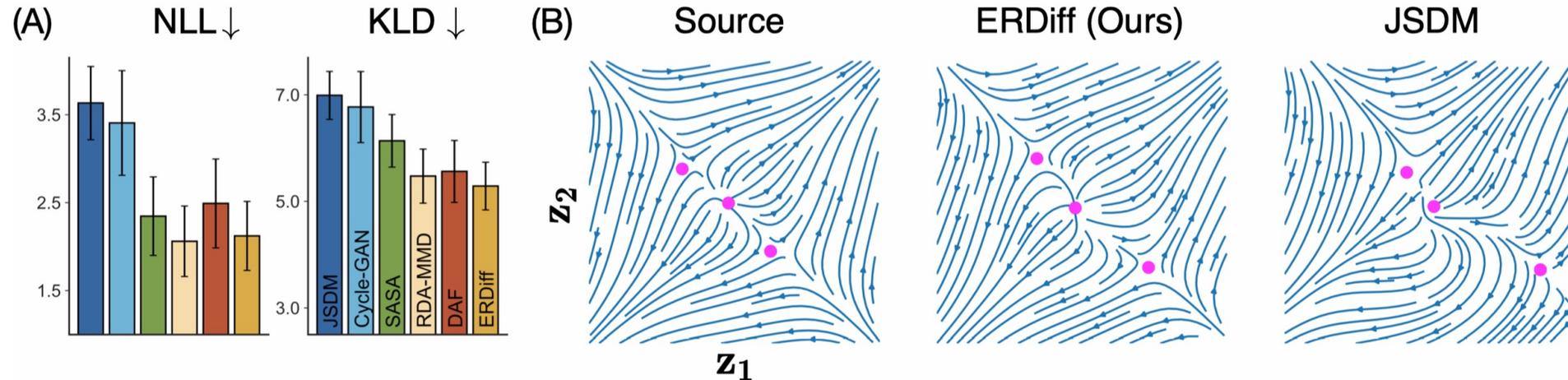
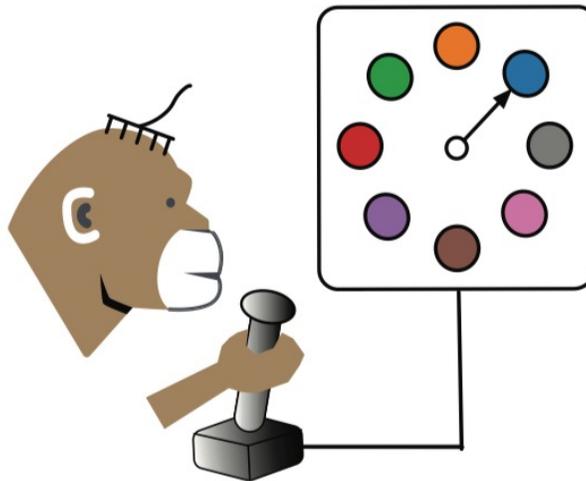


Figure 3: **Experimental Results on the synthetic dataset.** (A) Performance comparison on trial-average negative log-likelihood (NLL) and KL Divergence (KLD).  $\downarrow$  means the lower the better. ERDiff achieves the second-lowest NLL and the lowest KLD. (B) True continuous Bernoulli dynamics in the source domain compared to the latent dynamics aligned by ERDiff and JSMD in the target domain (blue dots denote the fixed points). ERDiff preserves the spatio-temporal structure of latent dynamics much better.

# Experiment

## Neural Datasets – Setup

- We conduct experiments with datasets collected from the primary motor cortex (M1) of non-human primates
- The primates have been trained to reach one of eight targets at different angles
- We perform cross-day (recordings of the same primate performing the task on different days) and inter-subject (recordings of different primates) experiments.



# Neural Datasets – Neural Manifold Analysis

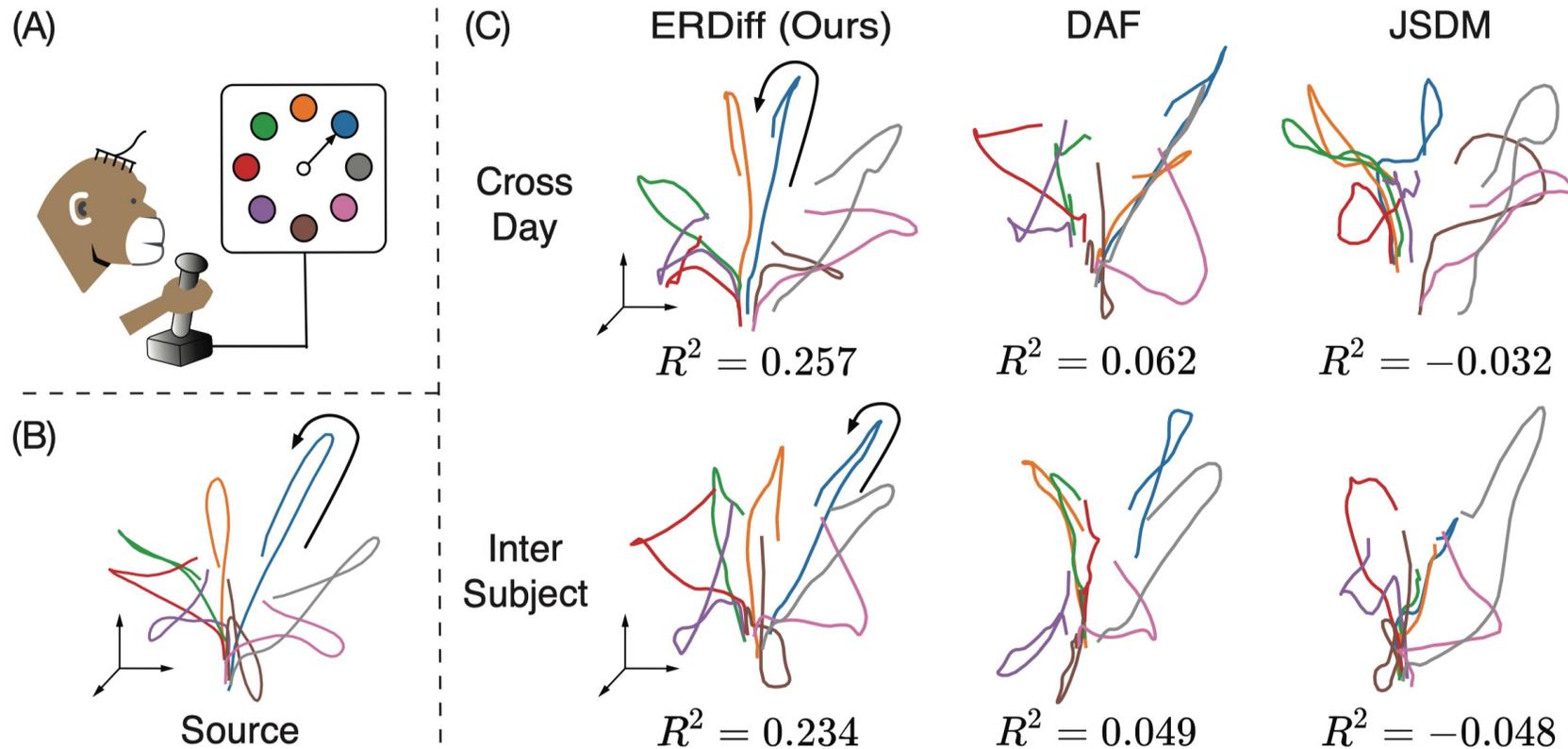


Figure 4: **Motor cortex dataset and Experimental Results.** (A) Illustration of the center-out reaching task of non-human primates. (B) The 3D Visualization of trial-averaged latent dynamics corresponding to each reaching direction in the source domain. (C) The 3D Visualization of trial-averaged latent dynamics corresponding to each reaching direction aligned by ERDiff, DAF, and JSMD given the target distribution from cross-day and inter-subject settings. We observe that ERDiff preserves the spatio-temporal structure of latent dynamics well.

# Neural Datasets – Behavior Decoding Performance

Table 1: The R-squared values ( $R^2$ , in %) and RMSE of the methods on the motor cortex dataset. ERDiff w/o S is short for a variant of our proposed method that removes the spatial transformer layer in the DM. ERDiff w/o T is short for a variant of our proposed method that removes the temporal transformer layer in the DM. The boldface denotes the highest score. Each experiment condition is repeated with 5 runs, and their mean and standard deviation are listed.

Method	Cross-Day		Inter-Subject	
	$R^2$ (%) $\uparrow$	RMSE $\downarrow$	$R^2$ (%) $\uparrow$	RMSE $\downarrow$
Cycle-GAN	-24.83 ( $\pm 3.91$ )	11.28 ( $\pm 0.44$ )	-25.47 ( $\pm 3.87$ )	12.23 ( $\pm 0.46$ )
JSDM	-17.36 ( $\pm 2.57$ )	9.01 ( $\pm 0.38$ )	-19.59 ( $\pm 2.77$ )	11.55 ( $\pm 0.52$ )
SASA	-12.66 ( $\pm 2.40$ )	8.36 ( $\pm 0.32$ )	-14.33 ( $\pm 3.05$ )	10.62 ( $\pm 0.40$ )
DANN	-12.57 ( $\pm 3.28$ )	8.28 ( $\pm 0.32$ )	-18.37 ( $\pm 3.24$ )	10.66 ( $\pm 0.57$ )
RDA-MMD	-9.96 ( $\pm 2.63$ )	8.51 ( $\pm 0.31$ )	-6.31 ( $\pm 2.19$ )	10.29 ( $\pm 0.42$ )
DAF	-6.37 ( $\pm 3.72$ )	8.17 ( $\pm 0.48$ )	-11.26 ( $\pm 3.64$ )	<b>9.57</b> ( $\pm 0.58$ )
ERDiff w/o S	-12.69 ( $\pm 2.64$ )	8.57 ( $\pm 0.50$ )	-14.60 ( $\pm 2.88$ )	10.85 ( $\pm 0.57$ )
ERDiff w/o T	-14.61 ( $\pm 2.33$ )	8.93 ( $\pm 0.50$ )	-17.10 ( $\pm 3.23$ )	10.94 ( $\pm 0.59$ )
<b>ERDiff (Ours)</b>	<b>18.81</b> ( $\pm 2.24$ )	<b>7.99</b> ( $\pm 0.43$ )	<b>10.29</b> ( $\pm 2.86$ )	9.78( $\pm 0.50$ )

# Computational Cost and Hyper-parameter Generalization

Table 2: Comparative analyses of computational cost between ERDiff and baseline methods during alignment. ERDiff has a comparable computational cost and maintains the stability of alignment.

Method	Cycle-GAN	JSDM	SASA	RDA-MMD	DAF	<b>ERDiff</b>
Add'l. Param	26K	0K	33K	65K	91K	28K
Add'l. Size	117KB	0KB	187KB	314KB	367KB	139KB
Align. Time	103ms	77ms	155ms	264ms	251ms	183ms
Stability	✗	✓	✓	✗	✗	✓

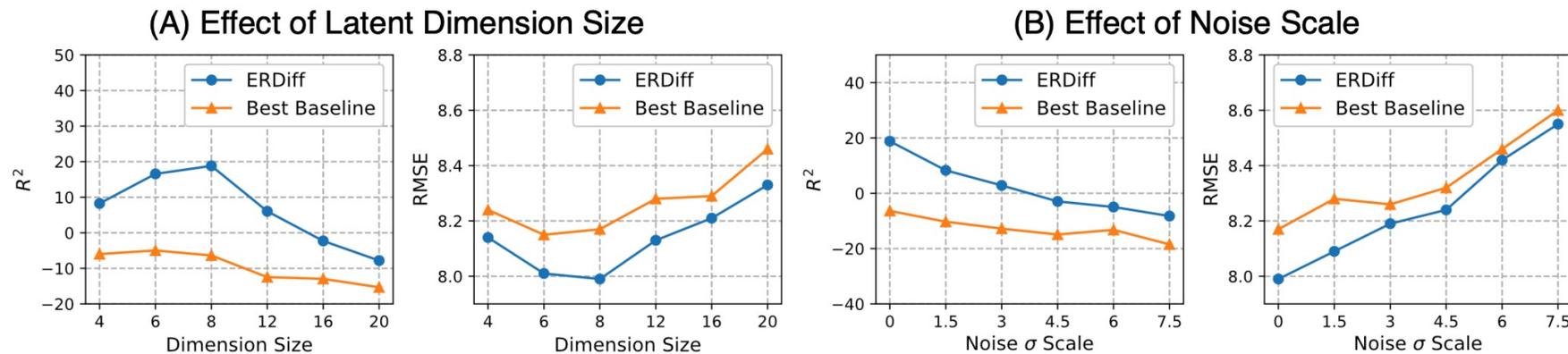


Figure 2: **Generalizability** and **Robustness** of ERDiff on alternative latent dimension sizes and noise scales (with 5 different random seeds). ERDiff consistently outperforms the best baseline method.

The background of the slide is a faded, light-colored photograph of a large, multi-story brick building with many windows, situated on a hillside. In the foreground, a paved path leads up the hill, and several people are walking away from the camera. The overall tone is warm and professional.

# Thanks for listening!