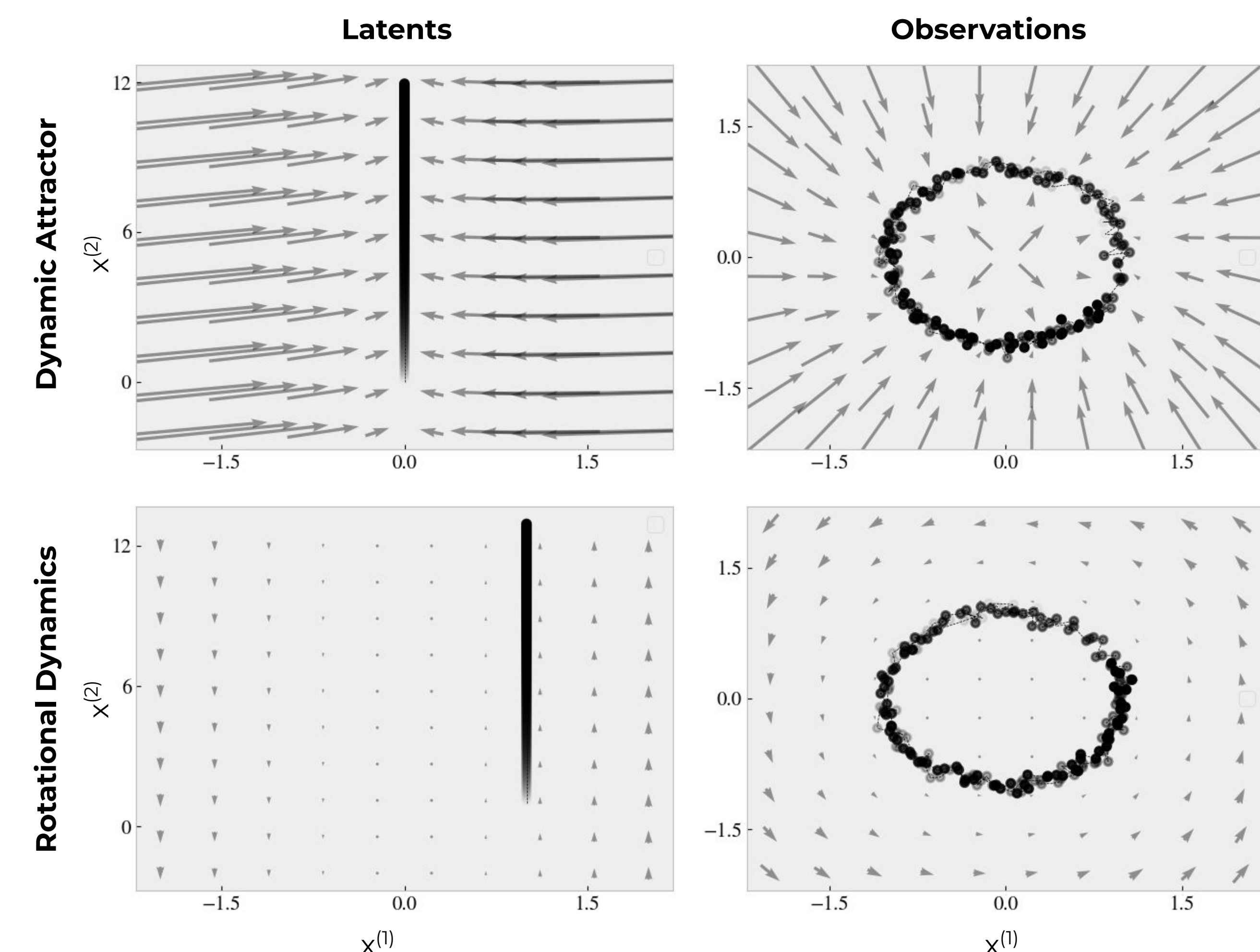


Introduction: Identifying Dynamics

Questions

- Which dynamical system model generated my data?
- Is motor cortex using continuous attractors or not?
- What are the latent variables underlying the observations?

The two models below are indistinguishable from the observational data alone



Theoretical Results: iSSM is Identifiable

Interventional State Space Models

- Latents** $\mathbf{x}_{t+1} = \mathbf{1}\{\mathbf{B}\mathbf{u}_t = 0\} \otimes \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \epsilon_t$ **Observations** $\mathbf{y}_t \sim P(\mathbf{y}_t | \mathbf{f}_\theta(\mathbf{x}_t))$.
- $\mathbf{u}_t \in \mathbb{R}^M$ interventional input to individual channels at time t .
 - $\mathbf{y}_t \in \mathbb{R}^N$ neural responses at time t , e.g. N -vector that concatenates the spike counts or calcium activities of all neurons.
 - $\mathbf{x}_t \in \mathbb{R}^D$ D -dimensional time-dependent latent variable.
 - $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ and \otimes denotes element-wise multiplication; $\mathbf{f}_\theta(\cdot)$ generic nonlinear function mapping latents to observations.
 - $\mathbf{A} \in \mathbb{R}^{D \times D}$ captures latent dynamics; $\mathbf{B} \in \mathbb{R}^{D \times M}$ captures the effect of neural perturbations on latent dynamics.
 - If the intervention \mathbf{u}_t is zero, the model follows observational dynamics.
 - In the presence of an intervention, the model decouples the intervened node from its parents.

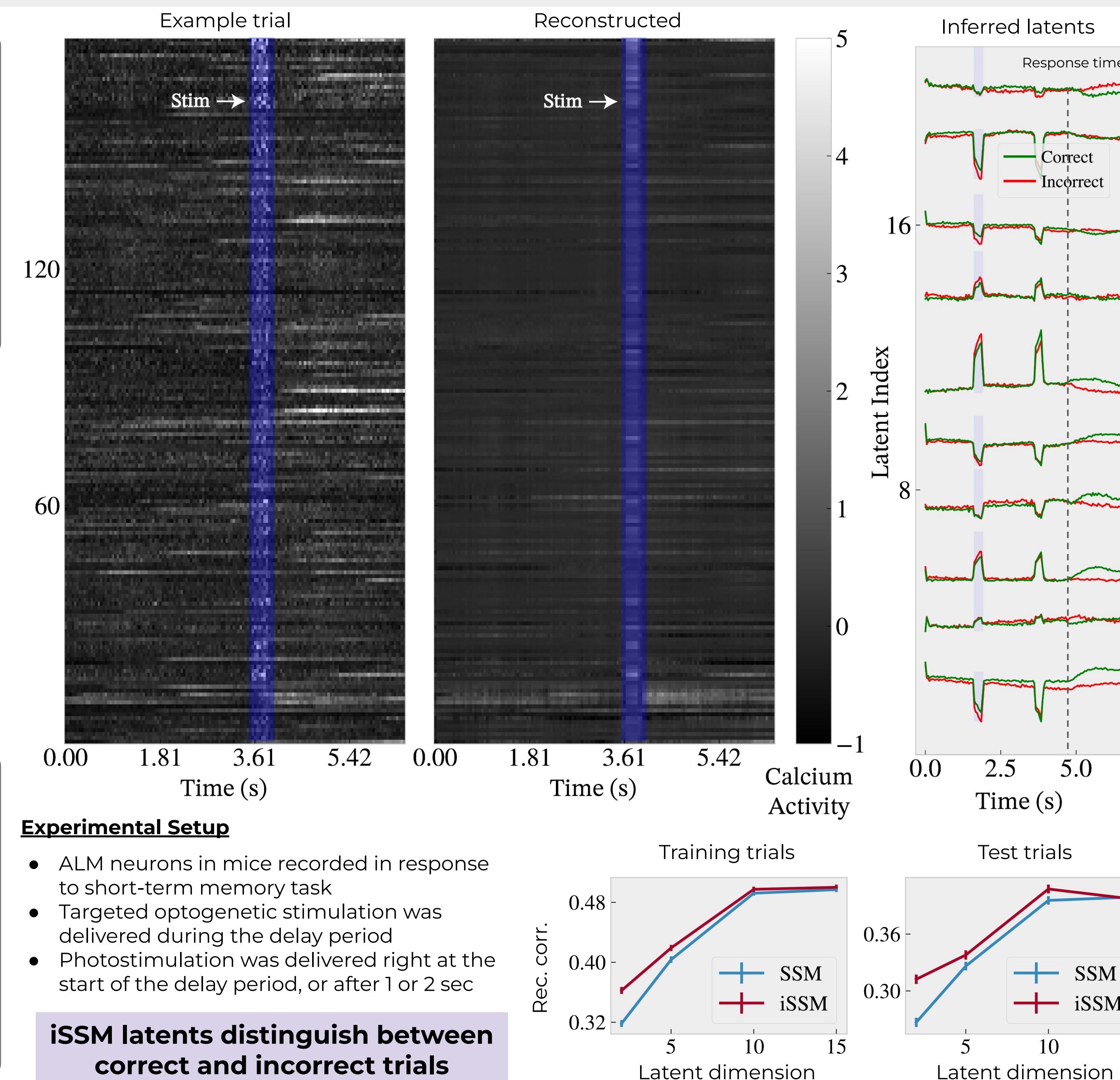
Identification Assumptions

- Assumption 1 (Observation model).** The function $P(\mathbf{y}_t | \mathbf{z}_t)$, where $\mathbf{z}_t = \mathbf{f}_\theta(\mathbf{x}_t)$, is **bounded complete** in \mathbf{y}_t .
 - [e.g. including exponential families, location-scale families, and nonparametric regression models]
- Assumption 2 (Mixing function).** The mixing function $\mathbf{f}_\theta(\cdot)$ is piecewise linear, continuous, and injective.
 - [e.g. including (deep) ReLU networks]
- Assumption 3 (Faithfulness).** There does not exist a non-zero vector \mathbf{V} such that $\text{Cov}(\mathbf{V}^\top \mathbf{x}_{t+1}, \mathbf{V}^\top \mathbf{x}_t) = 0, \forall t$
 - [Loosely, each latent dimension has at least one (non-trivial) causal parent from the previous timestep]

Identifiability Guarantees

- Theorem (Block identifiability of iSSM and generalization to unseen interventions).**
 - Under Assumptions 1-3, the latent dynamics \mathbf{A} and the mixing function of $\mathbf{f}_\theta(\cdot)$ can be **block-identified** up to permutation, and **shifting** and **scaling**.
 - Given a single intervention trial, one can **separate out** the intervened latents from the un-intervened ones.
 - Can extrapolate to **novel, unseen interventions** as long as they only touch upon already separated latents.
- Corollary (Identifiability of iSSM under sufficiently diverse interventions).**
 - If the interventions satisfy the **unordered pairs condition** (Hytinen et al., 2013), then the iSSM is **identifiable** up to permutation, along with coordinate-wise scaling and shifting.
 - The distribution under **any novel interventions** is also **identifiable**.

Results: Optogenetics in Mouse ALM



Results: Models of Working Memory

Experimental Setup

- Low-rank (LR) vs. functionally feed-forward (FF) are proposed as models of persistent activity in working memory
- Data was generated from LR and FF models in response to stimulation
- iSSM and SSM were fit using linear dynamics and linear observation models

Noisy linear dynamics

$$\frac{dx}{dt} = \mathbf{A}\mathbf{x} + d\mathbf{w}_t$$

LR connectivity

$$\mathbf{T} = \mathbf{I}_D^{(0)} + 0.5 [\mathbf{0}_{D \times (p-1)} \quad \mathbf{1}_{D \times 1}]$$

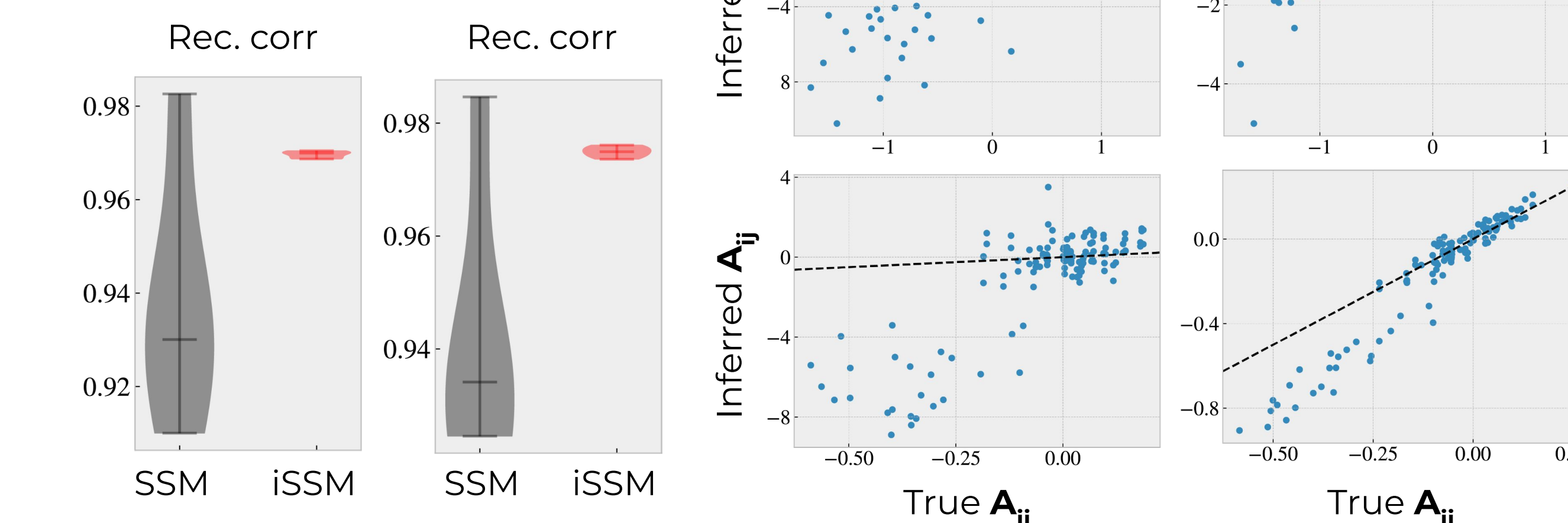
$$\mathbf{A} = -\mathbf{I}_D + \mathbf{O}\mathbf{T}\mathbf{O}^\top \quad \mathbf{O} \text{ is a randomly generated orthogonal matrix}$$

FF connectivity

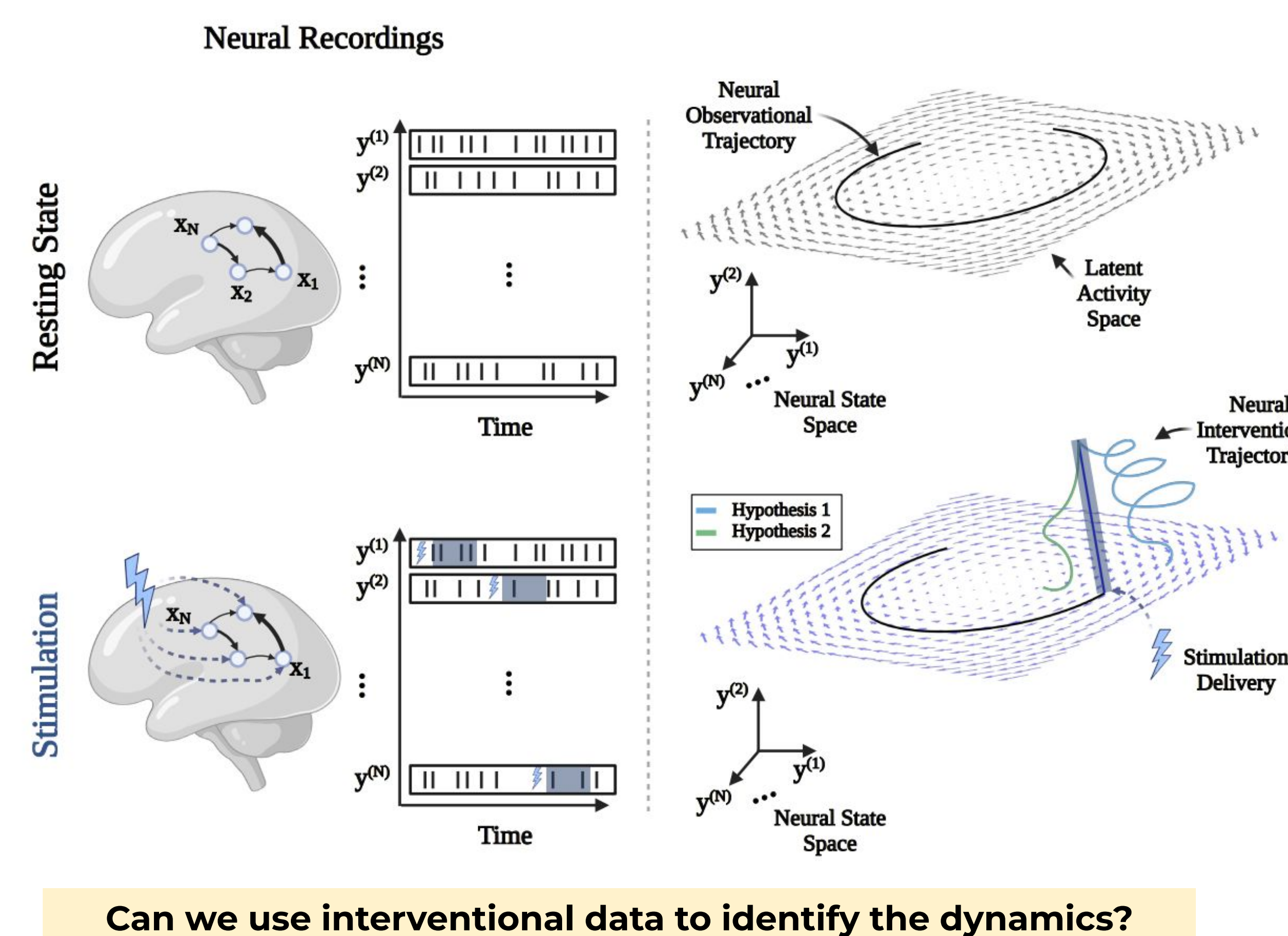
$$\mathbf{A} = \text{diag}(\exp(\ell_1), \exp(\ell_2), \dots, \exp(\ell_D)), \text{ where } \ell_i = \frac{D-i}{D-1}$$

$$\mathbf{A} = -\mathbf{I}_D + \mathbf{O}\mathbf{A}\mathbf{O}^\top$$

iSSM identifies the latents and the connectivity matrix



Approach: Interventional Models



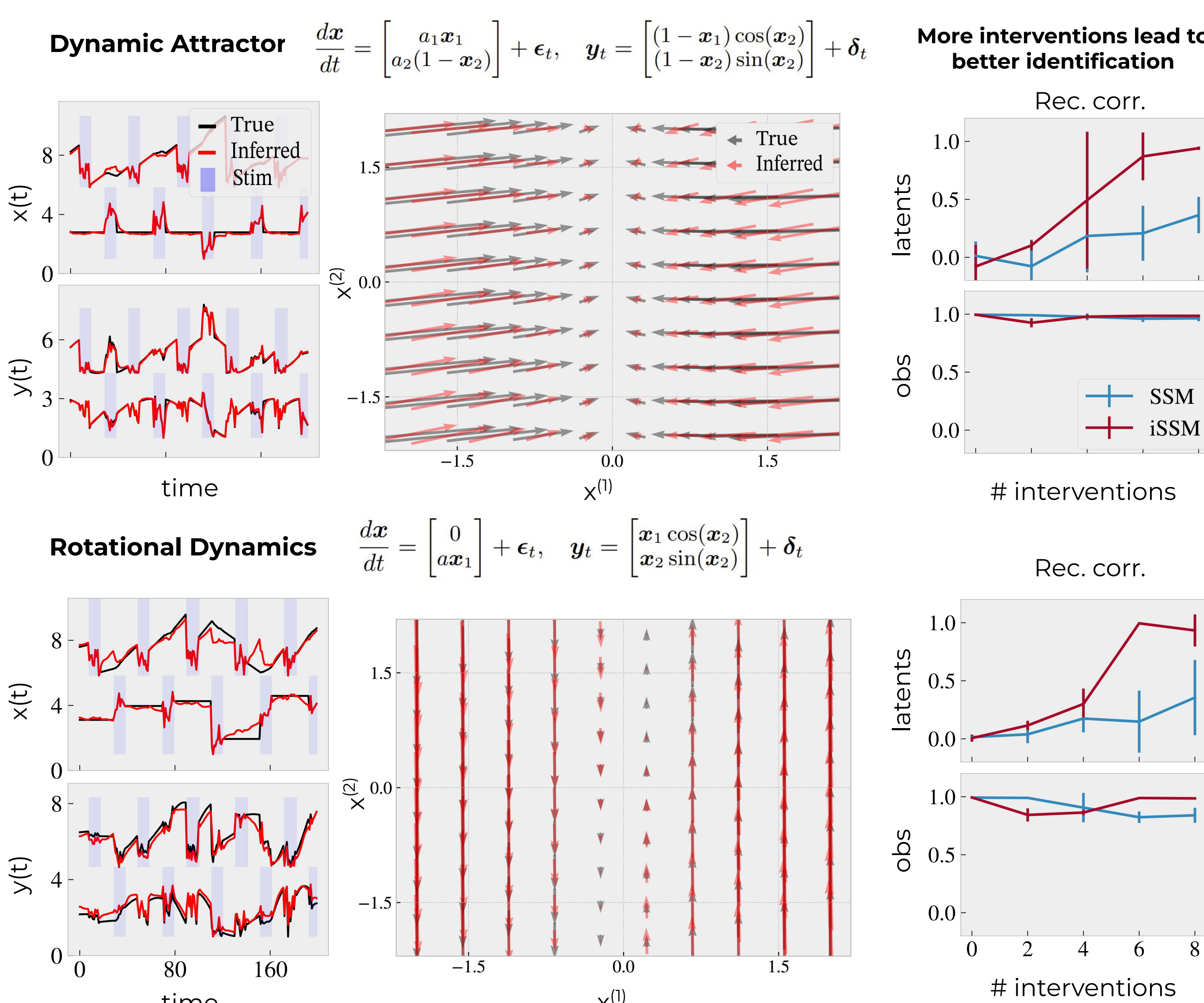
Can we use interventional data to identify the dynamics?

Intuition

- Interventions kick the state of the system outside of its attractor manifold, thereby allowing for the exploration of the state space and collecting more information about the dynamics
- However, interventional data alone is not sufficient for identification, we also need **interventional models** that properly leverage the interventional data

	Dynamics	Observations
SSM	$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \epsilon_t$	$\mathbf{y}_t \sim P(\mathbf{y}_t \mathbf{f}_\theta(\mathbf{x}_t))$
iSSM	$\mathbf{x}_{t+1} = \mathbf{1}\{\mathbf{B}\mathbf{u}_t = 0\} \otimes \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \epsilon_t$	$\mathbf{y}_t \sim P(\mathbf{y}_t \mathbf{f}_\theta(\mathbf{x}_t))$

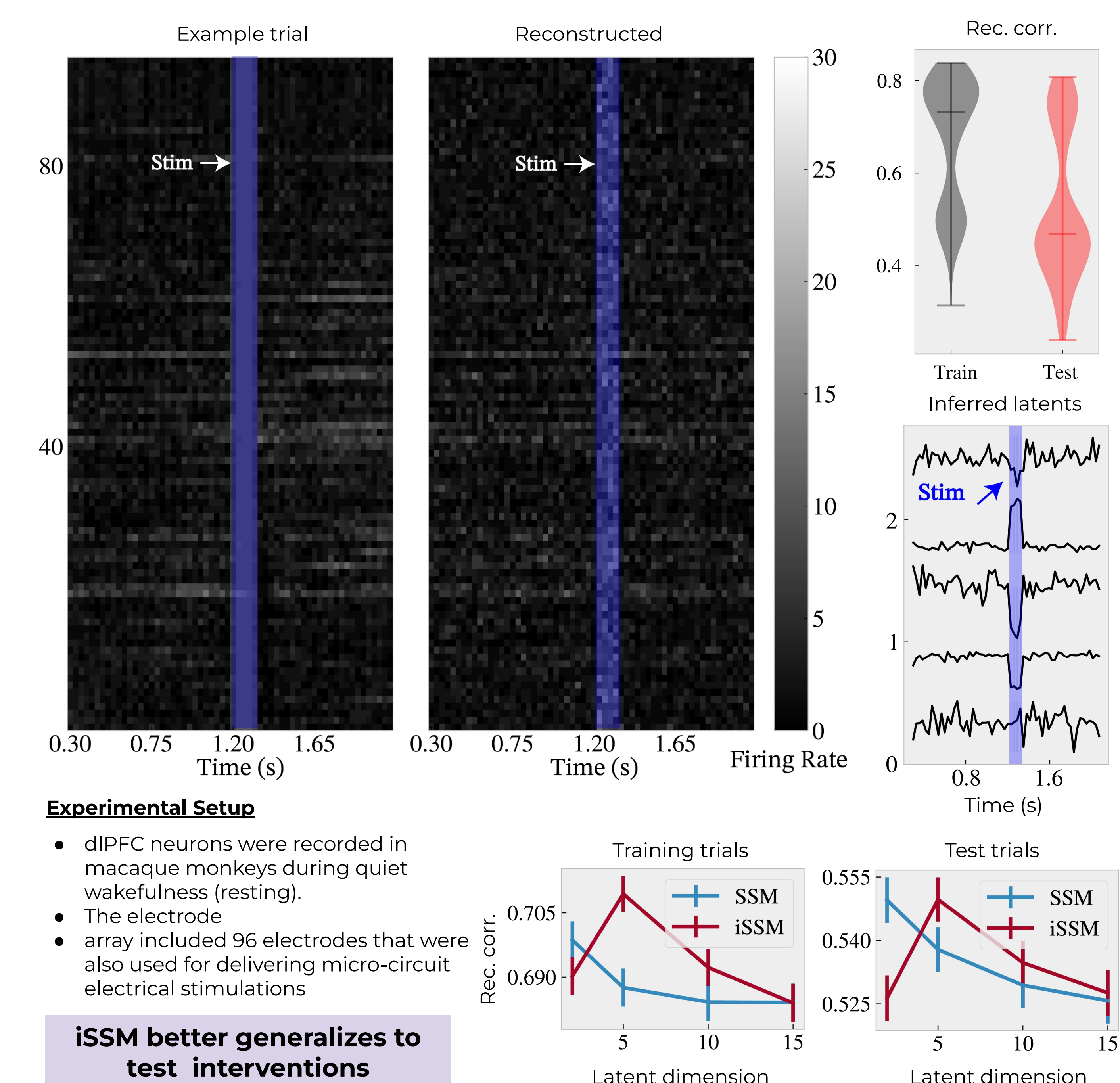
Results: Models of Motor Cortex



* For higher dimensional results check the paper

iSSM identifies the latents and the underlying flow field

Results: Micro-Stimulation in Primate dlPFC



Summary & References

Summary

- Here we proposed **iSSM**, a framework for joint modeling of observational and interventional data.
- We provided theoretical results showing that the iSSM model, when fitted on interventional data, leads to the **identifiability of latents as well as dynamics and emissions**.
- We showed in the **models of motor cortex** and **working memory** with **linear dynamics** and **linear or nonlinear emissions** iSSM leads to model identifiability.
- We showed an application of iSSM to **calcium recordings from the mouse ALM** region with **targeted photostimulation** delivered by channels that targeted groups of neurons.
- We showed an application of iSSM to **electrophysiological recordings** from the **macaque monkey prefrontal cortex with micro-stimulation** delivered by the same recording electrodes.

Future Directions

- (1) Interventional models with nonlinear dynamics
- (2) Modeling interventions applied to neurons as opposed to latents
- (3) Better inference algorithms

References

- Galgali, A., et al. Residual dynamics resolves recurrent contributions to neural computation. Nature Neuroscience, 26(2):326–338, 2023.
- Qian, W., et al. Partial observation can induce mechanistic mismatches in data-constrained models of neural dynamics. bioRxiv, pp. 2024–05, 2024.
- Daie, K., et al. Targeted photostimulation uncovers circuit motifs supporting short-term memory. Nature neuroscience, 24(2):259–265, 2021.
- Nejatbakhsh, A., et al. Predicting the effect of micro-stimulation on macaque prefrontal activity based on spontaneous circuit dynamics. Physical Review Research, 5(4):043211, 2023.

