

**Illustration:** DICA attempts to find representation  $s$  s.t.  $\partial \mathbf{x} / \partial s_i$  and  $\partial \mathbf{x} / \partial s_j$  spanning quite different directions, inducing larger convex hull volume (under constraints).

## Nonlinear Mixture Model Identification (NMMI) 🤖

A diffeomorphism mapping latent  $s$  to a  $d$ -dim. data manifold embedded in  $\mathbb{R}^m$ :

$$\mathbf{x} = \mathbf{f}(s), s \in \mathbb{R}^d, \mathbf{x} \in \mathbb{R}^m, d \leq m \quad (1)$$

- $s = [s_1, s_2, \dots, s_d] \sim p(s)$  are **latent** variables (object positions, lighting,...),
- $\mathbf{x} = [x_1, x_2, \dots, x_m]$  are **observed** features (pixels).
- $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^m$  is nonlinear mixing function.

### Goal: Recovering of $s$ and $f$ (up to acceptable ambiguities)

Learn an encoder  $\mathbf{g}_\phi(\mathbf{x}) = \hat{s}$  such that  $\hat{s}_i = \rho_i(s_{\pi(i)}), \forall i \in \{1, \dots, d\}$ , for a permutation  $\pi(\cdot)$  & an element-wise invertible  $\rho_i(\cdot)$ .

Applications: disentanglement, causal representation learning, self-supervised learning, etc.

### Challenge: Non-identifiability

An infinite number of  $(f, s)$  can satisfy  $\mathbf{x} = \mathbf{f}(s)$  in (1) 😞

## Related Works

Identifiability challenge is notorious in **nonlinear ICA** (nICA): even with statistically independent  $s_1, \dots, s_d$ , the model  $\mathbf{x} = \mathbf{f}(s)$  is non-identifiable [1].

**nICA with Auxiliary Variables  $u$ .** Side information (time frame labels, observation group indices, view indices, etc.) can help underpin identifiability of NMMI via [2]

$$p(s|u) = \prod_{i=1}^d p(s_i|u). \quad (2)$$

😞 Diverse auxiliary  $u$  not always available

**nICA with Structured  $f$ .** Conformal/local isometry/post-nonlinear/piecewise-affine  $f$ .

😞 Structured  $f$  holds in limited applications

**Structured Jacobian.**  $[\mathbf{J}_f(s)]_{i,j} = \partial x_i / \partial s_j$  describes how  $x_i$  is influenced by  $s_j$ .

1. **Independent Mechanism Analysis:** orthogonal columns of  $\mathbf{J}_f(s)$  [3].

😞 Lacking global identifiability

2. **Structural Sparsity:** a sparsity pattern on  $\mathbf{J}_f(s)$ , proposes to minimize  $\|\mathbf{J}_f(s)\|_1$  [4].

3. **Object-centric Learning:** a compositional  $f$  — a sparsely structured  $\mathbf{J}_f(s)$  where non-zero blocks corresponds to an object in image [5].

😞  $\mathbf{J}_f(s)$  is non-sparse in many settings

## Sufficiently Diverse Influence (SDI) Condition

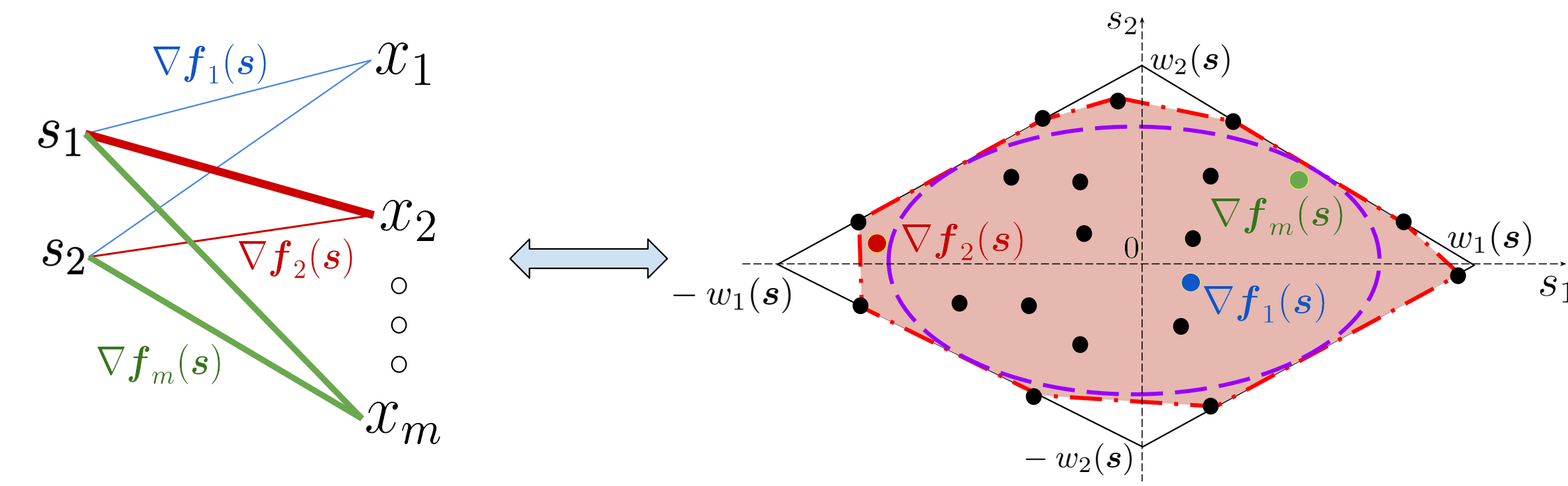
### Assumption: Sufficiently Diverse Influence

At  $s \in \mathcal{S}$ , there exists an  $s$ -dependent weighted  $L_1$ -norm ball  $\mathcal{B}_1^{w(s)}$  such that  $\nabla f_1(s), \dots, \nabla f_m(s) \in \mathcal{B}_1^{w(s)}$ . In addition:

1.  $\mathcal{E}(\mathcal{B}_1^{w(s)}) \subseteq \text{conv}\{\nabla f_1(s), \dots, \nabla f_m(s)\} \subseteq \mathcal{B}_1^{w(s)}$ ,
2.  $\text{conv}\{\nabla f_1(s), \dots, \nabla f_m(s)\}^* \cap \text{bd}(\mathcal{E}(\mathcal{B}_1^{w(s)}))^* = \text{extr}(\mathcal{B}_\infty^{w(s)})$ .

$\mathcal{S} \subset \mathbb{R}^d$ : set of latent variables;  $\mathcal{X} \subset \mathbb{R}^m$ : set of observations;  $\text{conv}\{\cdot\}$ : the convex hull of a set of vectors;  $\mathcal{E}(\mathcal{P})$  is its MVIE of polytope  $\mathcal{P}$ ;  $\mathcal{P}^*$ : polar set of  $\mathcal{P}$ ;  $\text{extr}(\mathcal{P})$  extreme points of  $\mathcal{P}$ ;  $\text{bd}(\mathcal{P})$ : boundary of  $\mathcal{P}$ .

## Illustration of Sufficiently Diverse Influence (SDI) Condition



*Visualizing Sufficiently Diverse Influence:* Row vectors  $\nabla f_1(s), \dots, \nabla f_m(s)$  of  $\mathbf{J}_f(s)$  are sufficiently distinct — their convex hull contains MVIE of an  $L_1$ -norm ball  $\mathcal{B}_1^{w(s)}$ .

**Origin.** SDI geometry originates from sufficiently-scattered condition (SSC) in NMF and PMF [6]; however, SSC characterizes the latent factors of a data matrix (e.g.,  $\mathbf{W}, \mathbf{H}$  in  $\mathbf{X} = \mathbf{W}\mathbf{H}$ ), do not involve nonlinear functions or derivatives as in SDI.

**Interpretation.** SDI reflects how  $s$  diversely affects  $x_1, \dots, x_m$ : some features are positively influenced by  $s_j$  ( $\partial x_i / \partial s_j > 0$ ), others are negatively influenced ( $\partial x_i / \partial s_j < 0$ ).

- Statistically dependent  $s$  and dense  $\mathbf{J}_f(s)$  can satisfy SDI.
- SDI favors  $m \gg d$  case (i.e., high-dim data with low-dim factors, say 📺/📺).
- SDI-satisfying geometric pattern of  $\nabla f_1(s), \dots, \nabla f_m(s)$  can vary with  $s$ .

## Learning Criterion: Jacobian Volume Maximization

Using  $\mathbf{f}_\theta, \mathbf{g}_\phi$  as two neural networks for autoencoder architecture  $\mathbf{x} = \mathbf{f}_\theta(\mathbf{g}_\phi(\mathbf{x}))$ .

$$\max_{\theta, \phi} \mathbb{E}[\log \det(\mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}))^\top \mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x})))] \quad (3)$$

$$\text{s.t. } \|\mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}))_{i,:}\|_1 \leq C, \quad \forall i = 1, \dots, m, \quad (4)$$

$$\mathbf{x} = \mathbf{f}_\theta(\mathbf{g}_\phi(\mathbf{x})), \quad \forall \mathbf{x} \in \mathcal{X} \quad (5)$$

- Objective (3) maximizes volume of convex hull of  $\mathbf{J}_{f_\theta}(\hat{s})$  spanned by its columns.
- Constraint (4) keeps the rows of  $\mathbf{J}_{f_\theta}(\hat{s})$  inside some  $L_1$ -norm ball, to comply with SDI.
- Constraint (5) keeps  $\mathbf{f}_\theta, \mathbf{g}_\phi$  invertible over  $d$ -dim manifold.

$\Rightarrow \partial \mathbf{x} / \partial s_1, \dots, \partial \mathbf{x} / \partial s_d$  are encouraged to scatter in space (inside a certain  $L_1$ -norm ball).

## Identifiability Results 🎯

### Identifiability of DICA

Let an optimal solution be  $(\hat{\theta}, \hat{\phi})$ . Assume  $\tilde{\mathbf{f}} = \mathbf{f}_{\hat{\theta}}$  and  $\tilde{\mathbf{g}} = \mathbf{g}_{\hat{\phi}}$  are universal function representers. Suppose the SDI condition holds for the NMMI model, for every  $s \in \mathcal{S}$ . Then,  $\tilde{s} = \tilde{\mathbf{g}}(\mathbf{x}) = \tilde{\mathbf{g}} \circ \mathbf{f}(s)$  where

$$[\tilde{s}]_i = [\tilde{\mathbf{g}}(\mathbf{x})]_i = \rho_i(s_{\pi(i)}), \quad \forall i \in [d], \quad (6)$$

in which  $\pi$  is a permutation of  $\{1, \dots, d\}$  and  $\rho_i(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is an invertible function.

### (Informal) Identifiability w/ Finite Number of SDI-Satisfying Points

Suppose each point in the finite set  $\mathcal{S}_N := \{s^{(1)}, \dots, s^{(N)}\}$  with  $\mathcal{X}_N := \{\mathbf{x} \in \mathcal{X} : \mathbf{x} = \mathbf{f}(s), \forall s \in \mathcal{S}_N\}$  is SDI-satisfying. Under several regularity conditions,  $\tilde{\mathbf{g}}(\mathbf{x}^{(n)}) = \tilde{\Pi} \tilde{\rho}(s^{(n)}), \forall n \in [N]$  for a constant permutation  $\tilde{\Pi}$ . With probability at least  $1 - \delta$ ,

$$\mathbb{E}_{s \sim p(s)} [\|\tilde{\mathbf{g}}(\mathbf{x}) - \tilde{\Pi} \tilde{\rho}(s)\|_2] = \mathcal{O}((L_f L_{\tilde{g}} + L_{\tilde{\rho}}) \mathcal{R}_N(\mathcal{G}) + \sqrt{\ln(1/\delta)/N}), \quad (7)$$

where  $\mathcal{R}_N(\mathcal{G})$  is the empirical Rademacher complexity of the encoder class.

## Implementation 🤖

Given  $L$  realizations of  $\mathbf{x}$ ,  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(L)}\}$ , use MLPs for  $\mathbf{f}_\theta, \mathbf{g}_\phi$ . At  $t$ -th epoch, optimize

$$\min_{\theta, \phi} \mathcal{L}_t := (1/L) \sum_{n=1}^L (\|\mathbf{x}^{(n)} - \mathbf{f}_\theta(\mathbf{g}_\phi(\mathbf{x}^{(n)}))\|_2^2 - \lambda_{\text{vol}}(t) \times c_{\text{vol}} + \lambda_{\text{norm}} \times c_{\text{norm}}(t)) \quad (8)$$

using a warm-up heuristic with  $T_w$  warm-up epochs:

- $c_{\text{vol}} := \log \det(\mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}^{(n)}))^\top \mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}^{(n)})))$  with  $\lambda_{\text{vol}}(t) := \frac{\lambda_{\text{vol}}}{T_w} \min\{t, T_w\}$  (a more computationally friendly trace-based surrogate of  $c_{\text{vol}}$  is available)

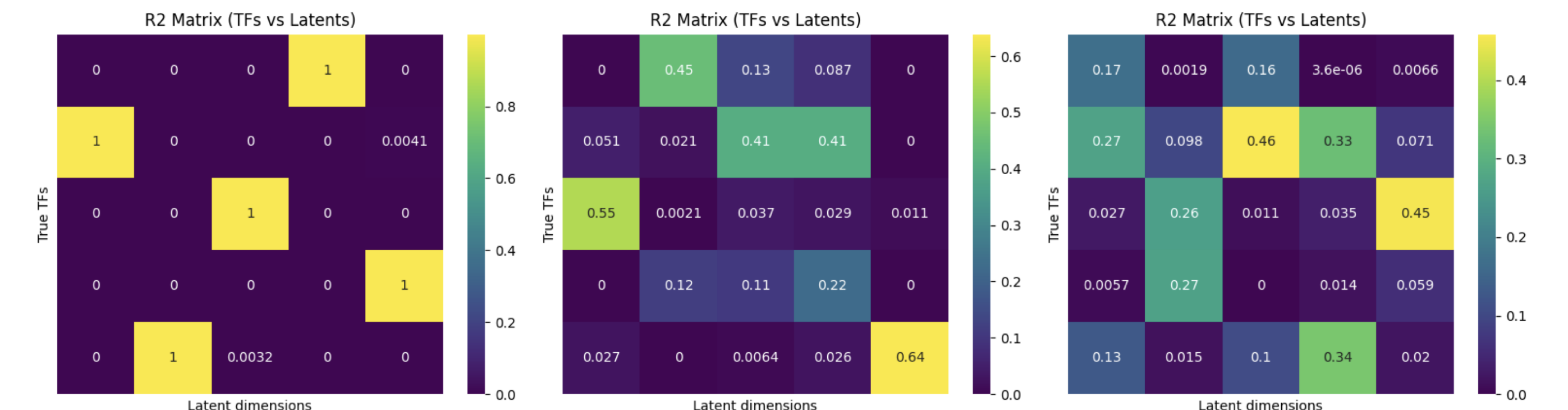
- $c_{\text{norm}}(t) := \begin{cases} \|\mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}^{(n)}))\|_1 & \text{if } t \leq T_w \\ \text{Softplus}\{\|\mathbf{J}_{f_\theta}(\mathbf{g}_\phi(\mathbf{x}^{(n)}))\|_1 - C\} & \text{if } t > T_w \end{cases}$  with  $\lambda_{\text{norm}} > 0$ , where  $C$  is average of last 10 epochs during warm-up.

## Experiments (more in our paper) 📊

### Single-cell Genomics 🧬

Inferring transcription factors' activities from gene expressions.

Data generation follows SERGIO simulator, using TRRUST dataset + cross-talk noises.

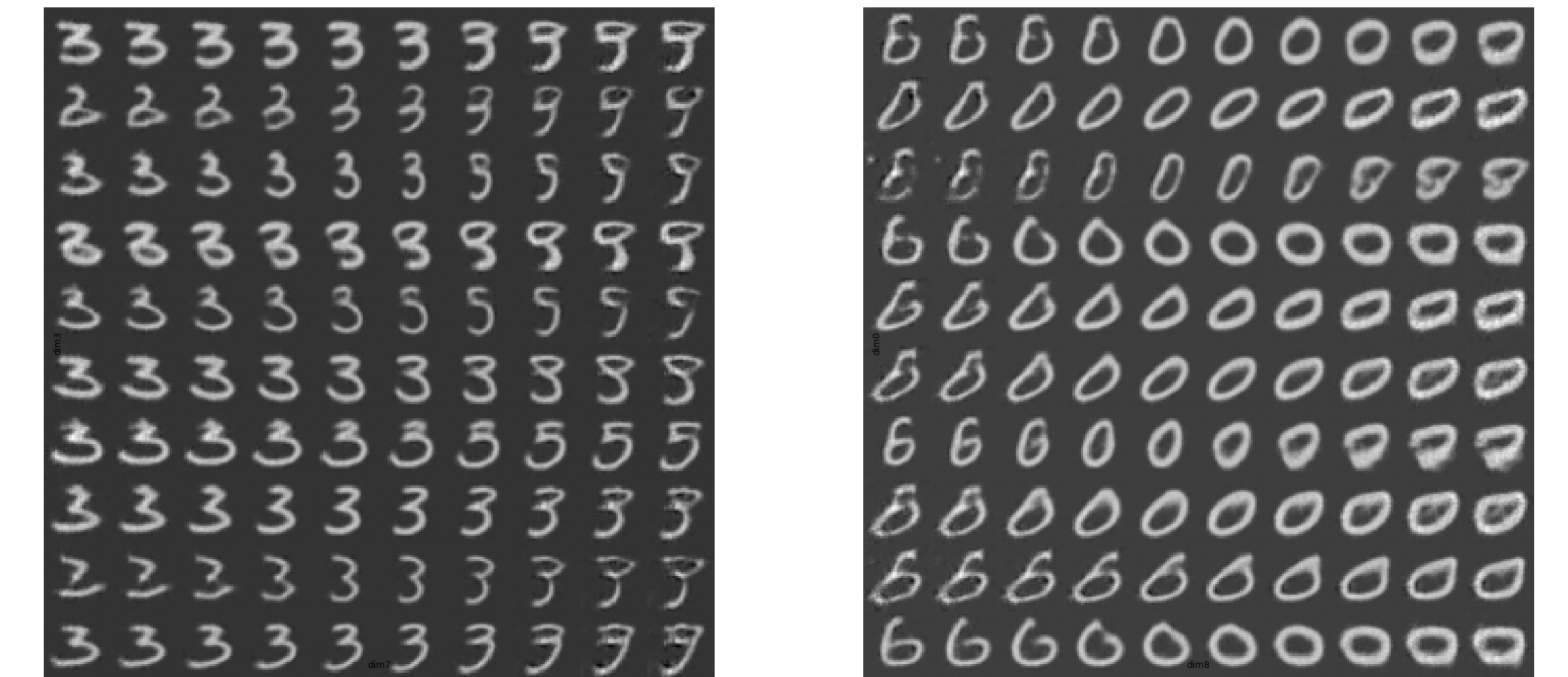


(a) DICA (mean  $R^2 \approx 1$ ) (b) Sparse (mean  $R^2 \approx 0.4544$ ) (c) Base (mean  $R^2 \approx 0.3381$ )

Heatmap of  $R^2$  scores between estimated components and ground-truth mRNA concentrations of TFs.

### Image Disentanglement 📺

Applying autoencoder with DICA loss on MNIST dataset.



(a) [Anchor digit 3] As  $s_8$  increases, digit 3 increasingly looks like digit 9. (b) [Anchor digit 0] As  $s_9$  decreases, digit 0 increasingly looks like digit 6.

### References:

- [1] Hyvärinen & Pajunen, "Nonlinear independent component analysis: Existence and uniqueness results", Neural Networks, 1999.
- [2] Hyvärinen et al., "Nonlinear ICA Using Auxiliary Variables and Generalized Contrastive Learning", AISTATS, 2019.
- [3] Gresle et al., "Independent mechanism analysis, a new concept?", NeurIPS, 2021.
- [4] Zheng & Zhang, "Generalizing Nonlinear ICA Beyond Structural Sparsity", NeurIPS, 2023.
- [5] Brady et al., "Provably Learning Object-Centric Representations", ICML, 2023.
- [6] Tatli & Erdogan, "Polytopic Matrix Factorization: Determinant Maximization Based Criterion and Identifiability", IEEE TSP, 2021.

**Acknowledgements:** This work is supported in part by the NSF CAREER Award ECCS-2144889.