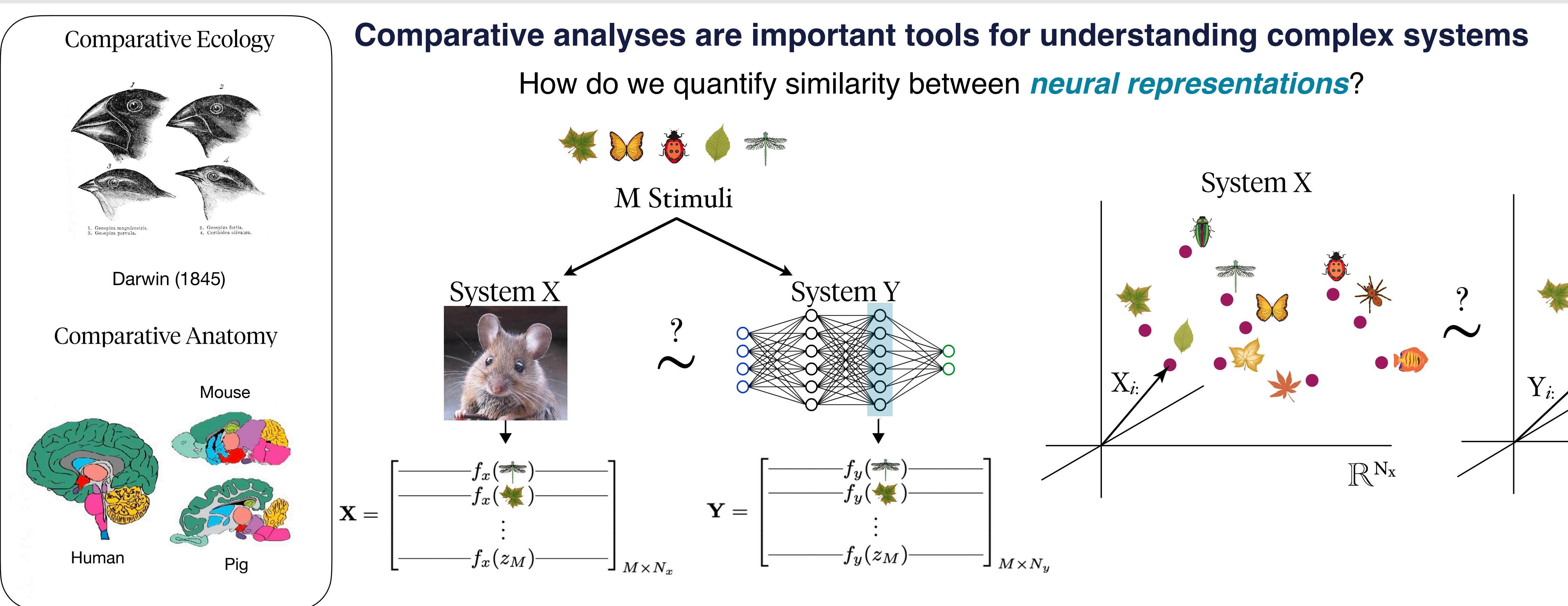
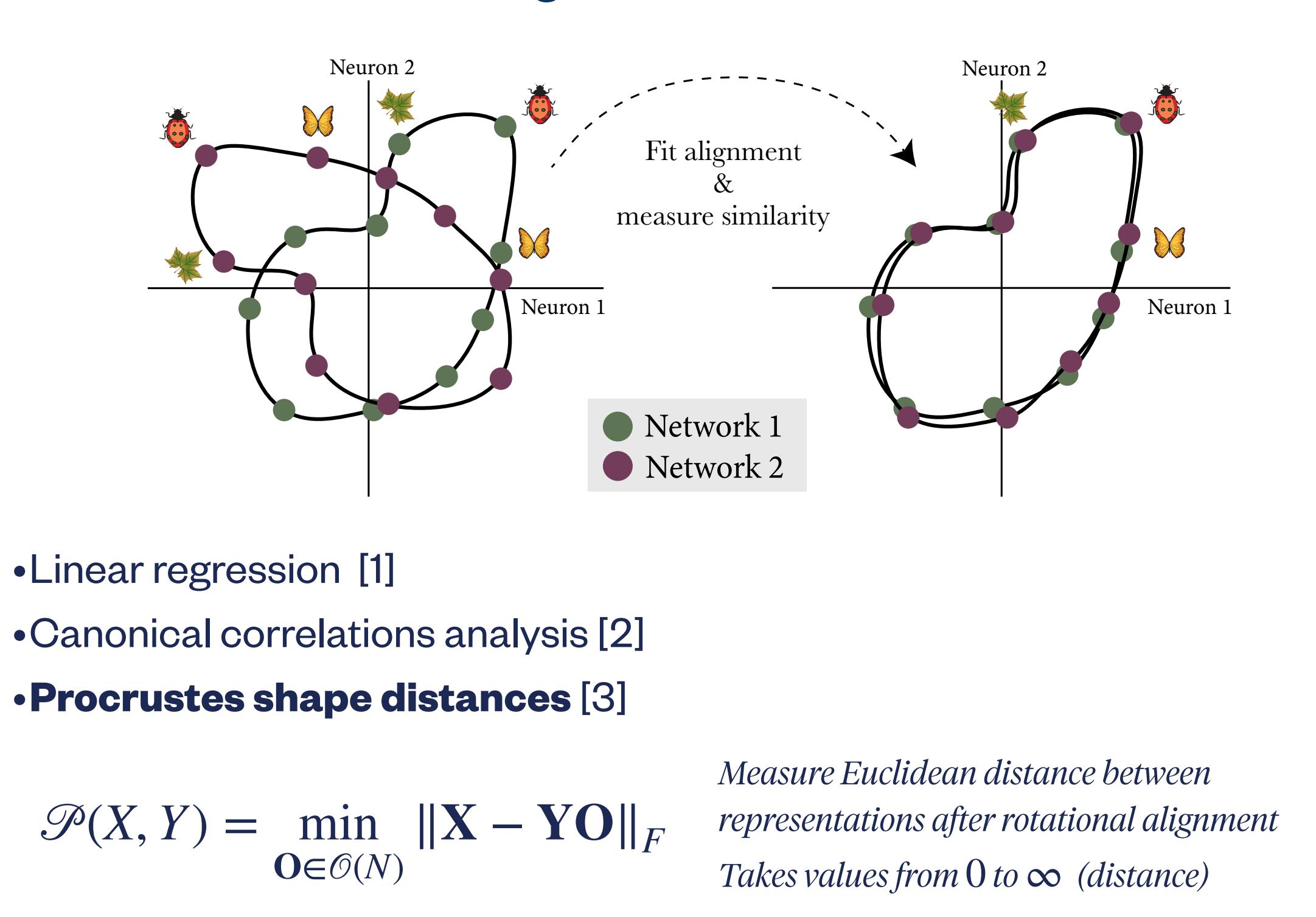


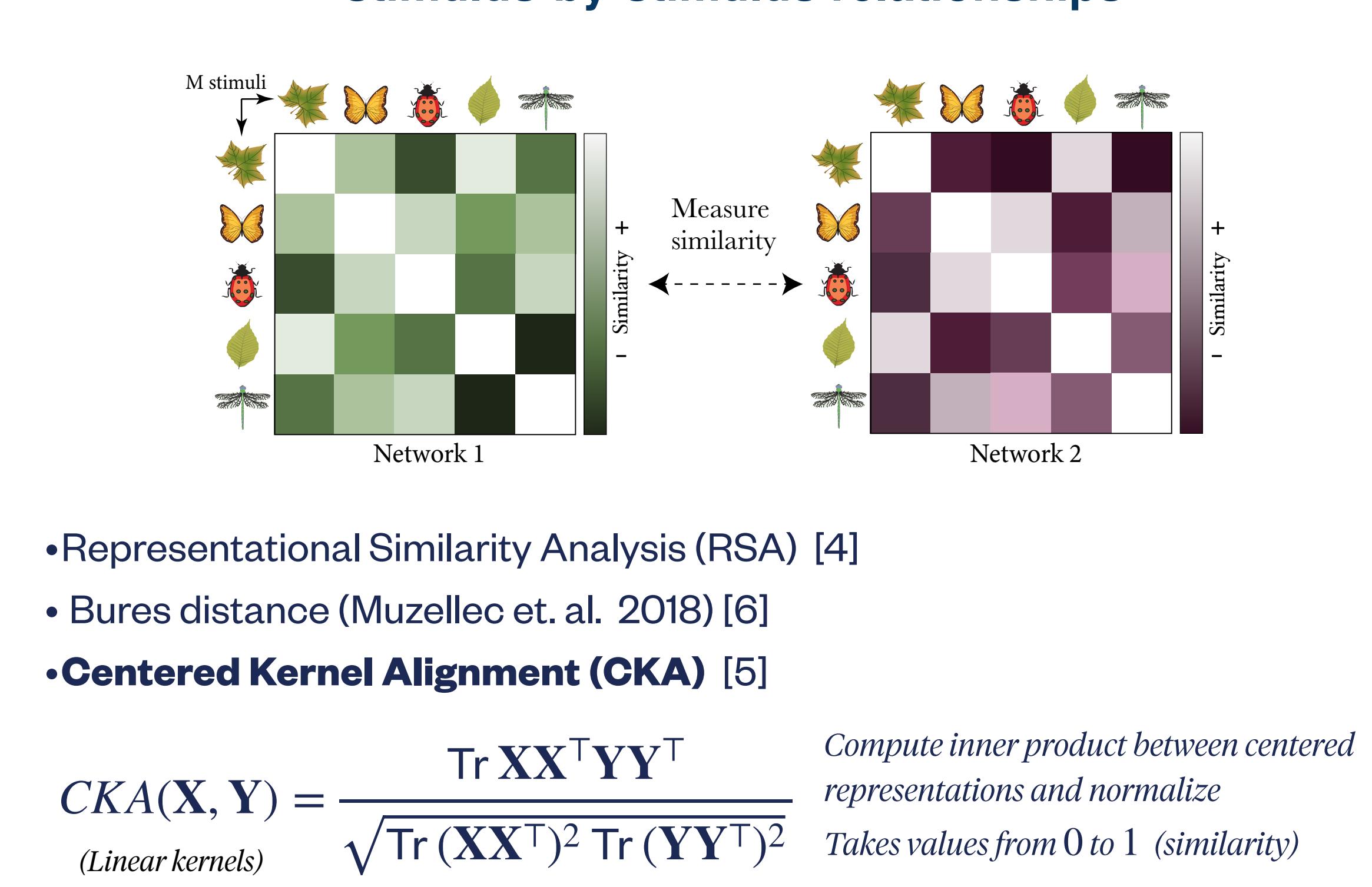
Measuring similarity between neural systems



(Dis)similarity measures that transform or align neural dimensions

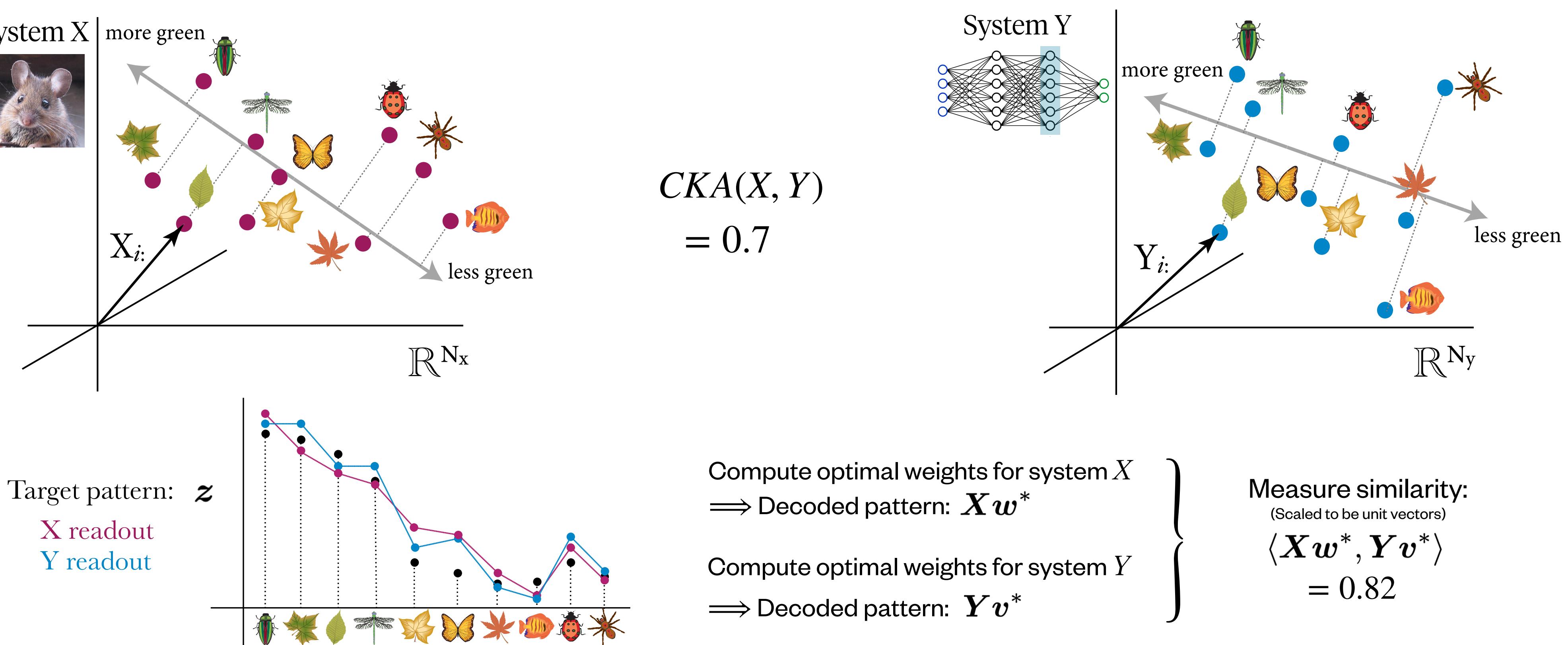


(Dis)similarity measures that quantify stimulus-by-stimulus relationships

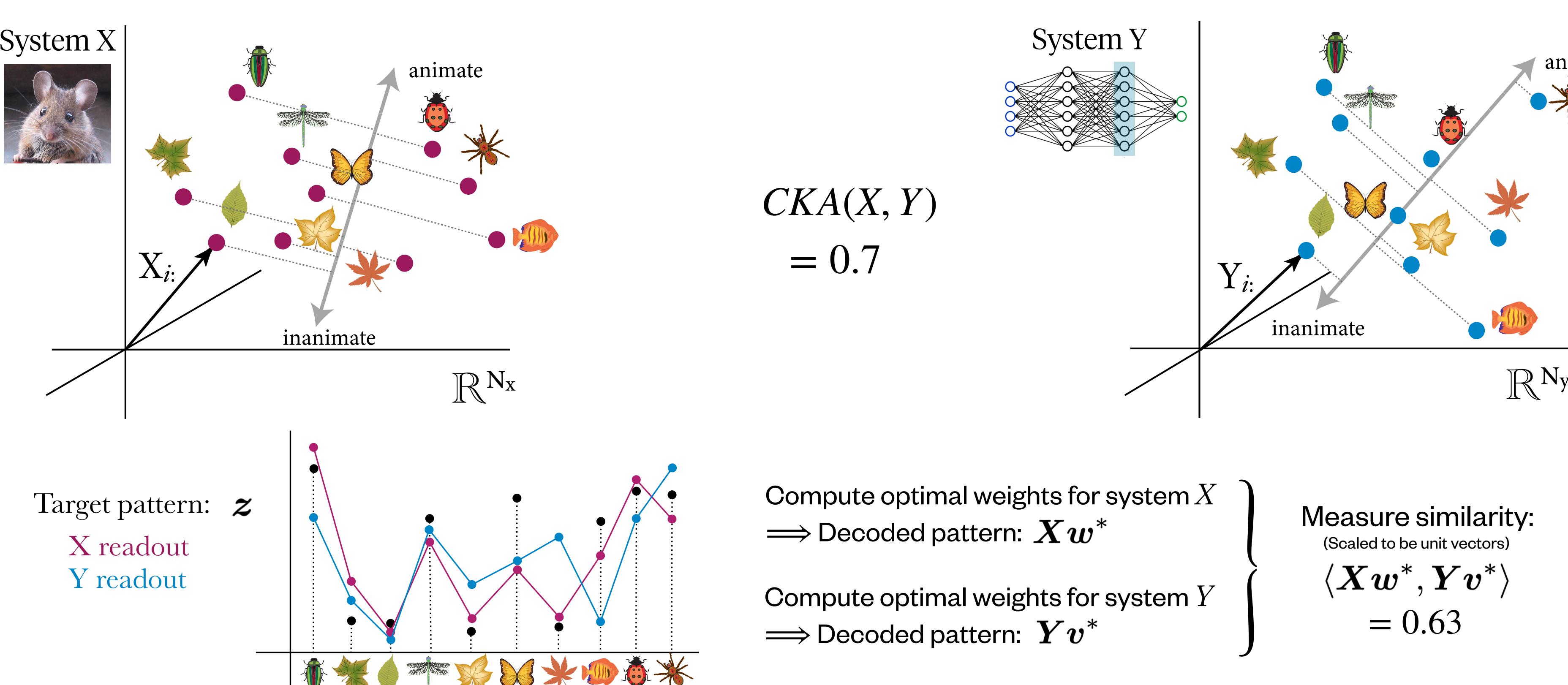


Comparing representations with linear decoding

Example task 1: color



Example task 2: animate/inanimate



Average decoding similarity/distance

★ Similarity depends on the choice of decoding task

$$\text{average decoding similarity (ADS)} \quad \mathbb{E}_{z \sim P_z} \langle Xw^*, Yv^* \rangle$$

To compute these, we must choose:

1. Regression loss function
2. Ensemble of tasks to average over

Idea: Measure similarity over an ensemble of decoding tasks

$$\text{average decoding distance (ADD)} \quad \mathbb{E}_{z \sim P_z} \|Xw^* - Yv^*\|_2^2$$

We show: certain choices here \Rightarrow average decoding similarity/distance = popular representational similarity/distance measures

Set up a family of linear decoding problems

Decoding optimization problem:

$$\underset{\mathbf{w}}{\text{maximize}} \quad \underbrace{\frac{1}{M} \mathbf{z}^T \mathbf{Xw}}_{\text{Maximize overlap between } \mathbf{Xw} \text{ and } \mathbf{z}} - \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{G}(\mathbf{X}) \mathbf{w}}_{\text{Penalty on a norm of the weights } \mathbf{w}}$$

This problem has a nice closed form solution:

$$\mathbf{w}^* = \frac{1}{M} \mathbf{G}(\mathbf{X})^{-1} \mathbf{X}^T \mathbf{z} \quad \text{Optimal Decoding Weights}$$

$\mathbf{G}(\cdot)$ is a function mapping $\mathbb{R}^{M \times N} \rightarrow$ symmetric positive definite $N \times N$ matrices

Consider $\mathbf{G}(\mathbf{X}) = a \mathbf{C}_X + b \mathbf{I}$ with neuron-by-neuron covariance $\mathbf{C}_X := \frac{1}{M} \mathbf{X}^T \mathbf{X}$

Relations to geometric similarity measures

Special cases

$$\text{Take } a = 1, b = \lambda \quad \mathbf{w}^* = \text{argmin} \|\mathbf{Xw} - \mathbf{z}\|_2^2 + \lambda \mathbf{I} \quad \text{Ridge regression}$$

$$\mathbf{v}^* = \text{argmin} \|\mathbf{Yw} - \mathbf{z}\|_2^2 + \lambda \mathbf{I}$$

If we also assume identity covariance structure of the task ensemble we are averaging over, i.e. $\mathbb{E}_{\mathbf{z}} [\mathbf{zz}^T] = \mathbf{I}$

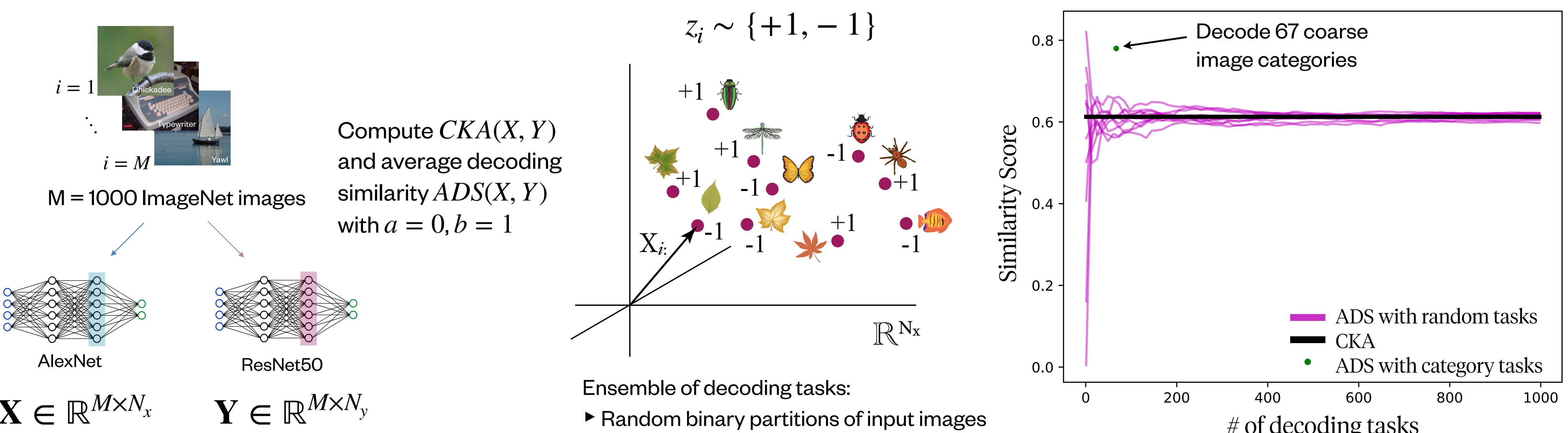
$$\text{Take } a = 0, b = 1 \quad \frac{\mathbb{E} \langle \mathbf{Xw}^*, \mathbf{Yv}^* \rangle}{\sqrt{\mathbb{E} \langle \mathbf{Xw}^*, \mathbf{Xw}^* \rangle \mathbb{E} \langle \mathbf{Yv}^*, \mathbf{Yv}^* \rangle}} = \text{CKA}(\mathbf{X}, \mathbf{Y})$$

(Normalized) Average decoding similarity

More in paper:

Similarity measure	a	b
Linear CKA	0	b
GULP	1	λ
CCA	1	0
ENSD	0	$\frac{1}{M} \text{Tr}[\mathbf{C}_X^2]$

Empirical Example: CKA

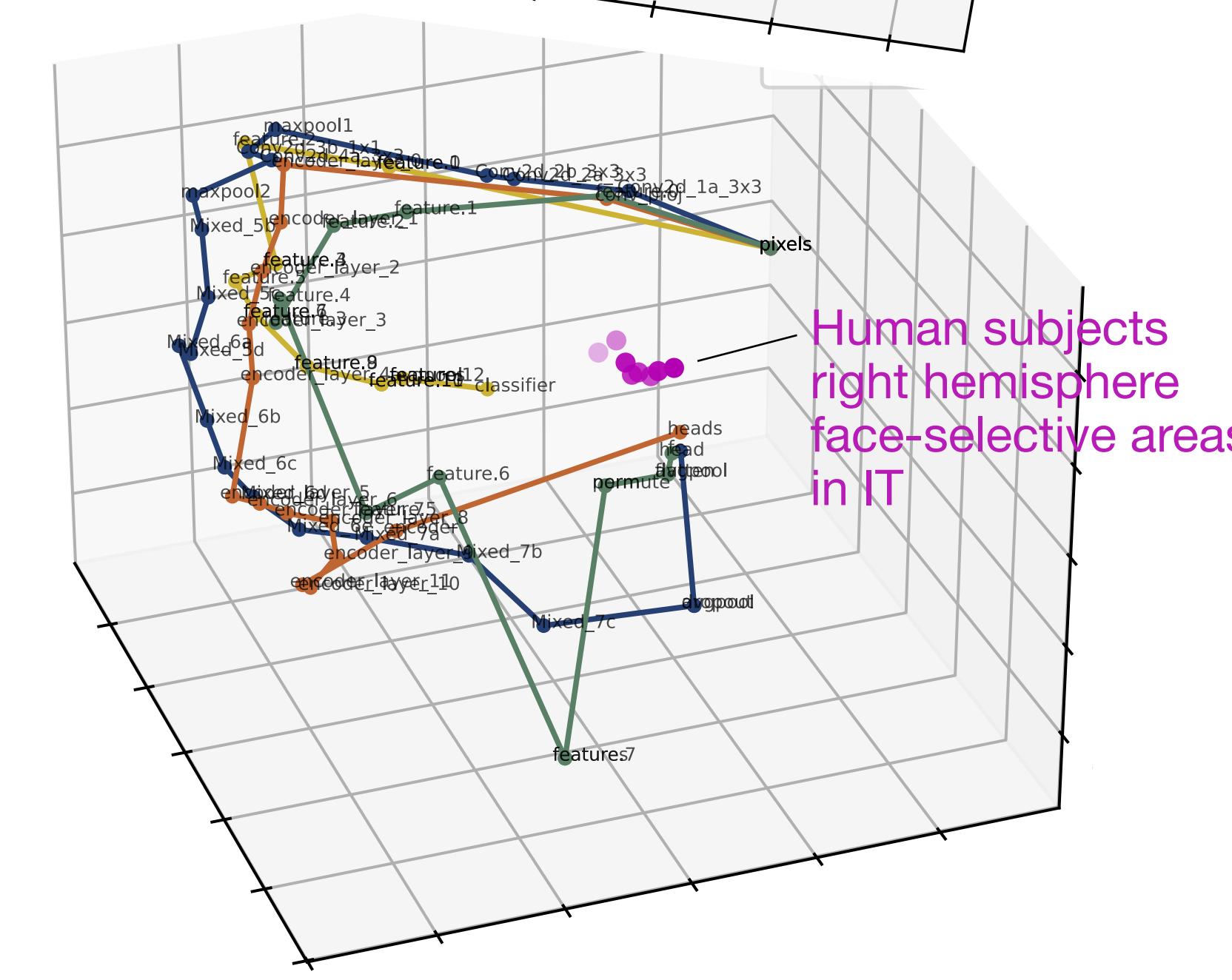
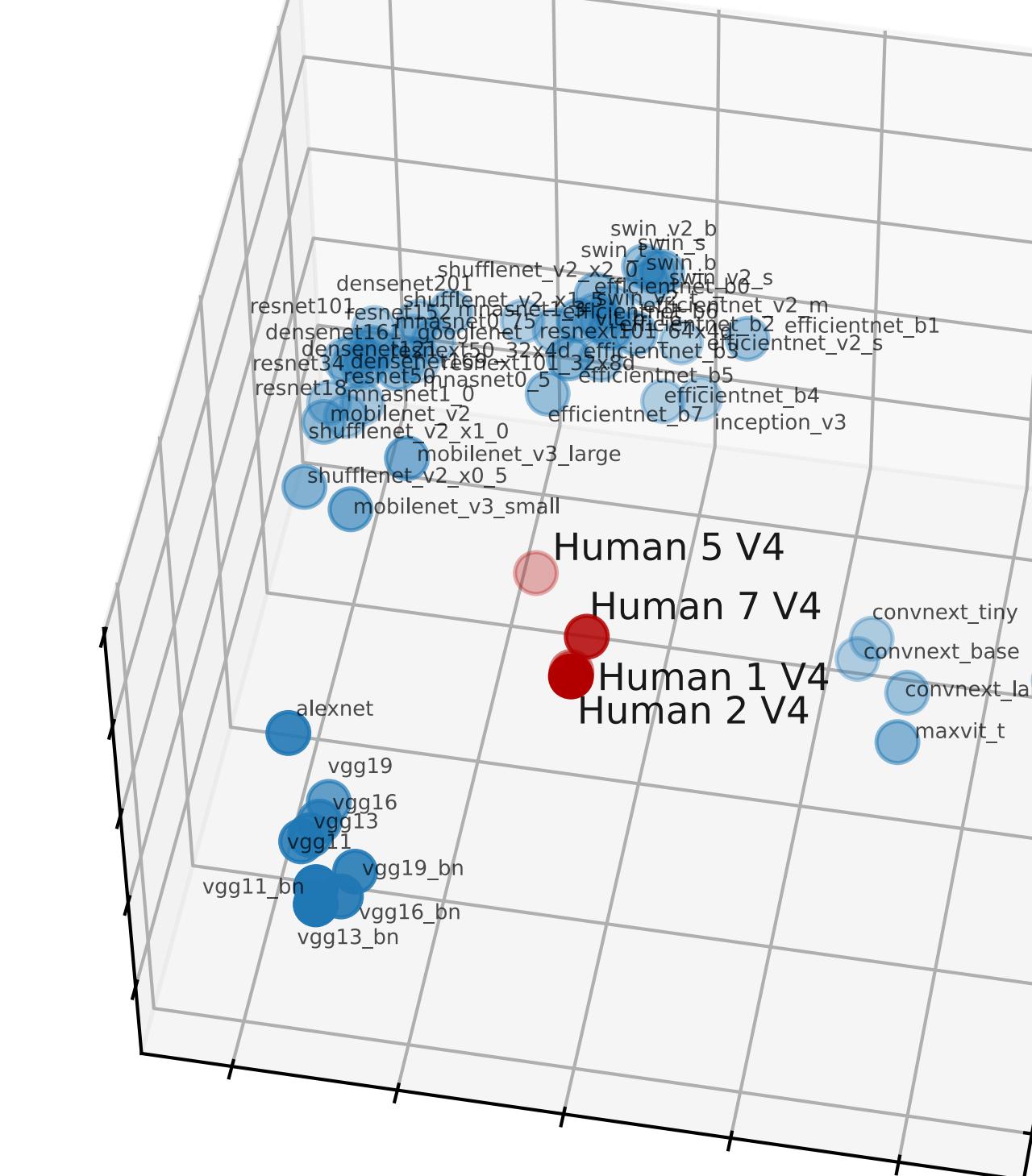
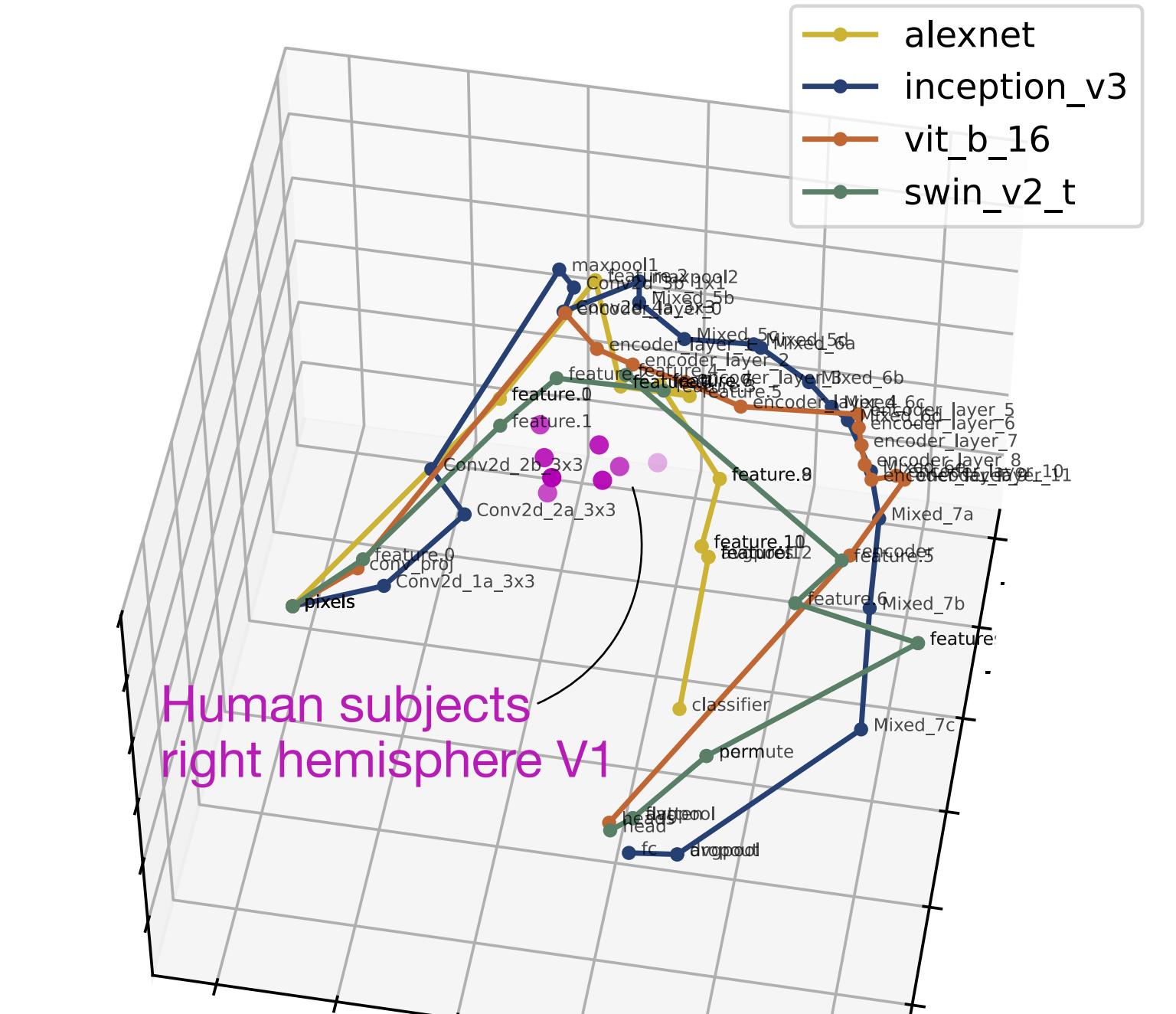
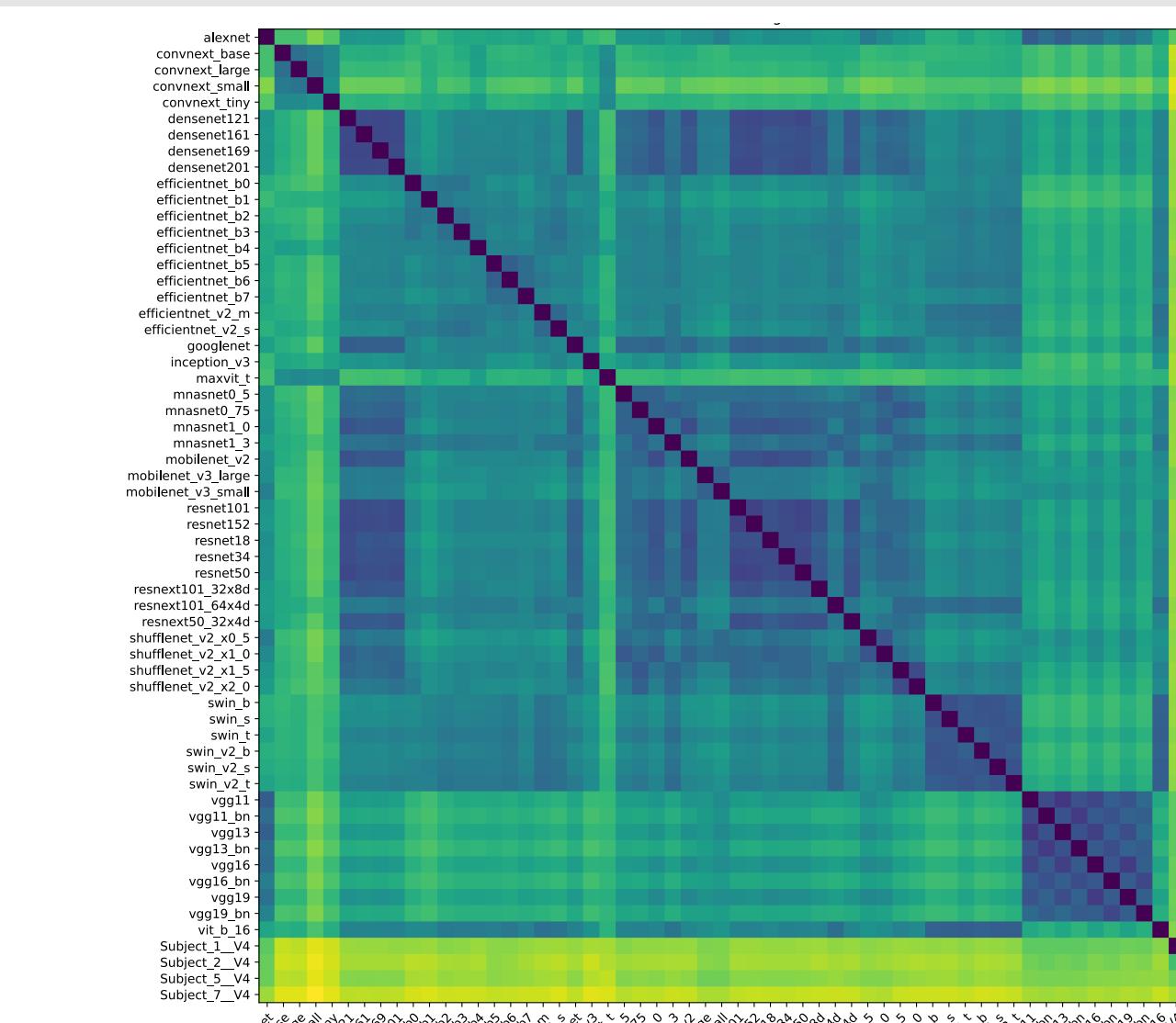


Example: Comparing deep network representations and human fMRI data

Average Decoding Similarity with human target

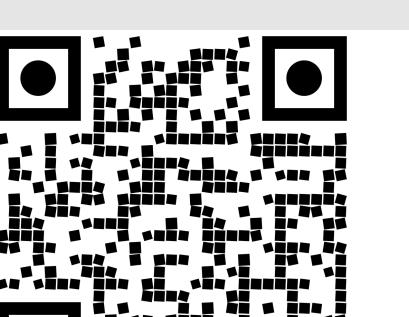
$$\mathbb{E}_{\mathbf{z}} [\mathbf{zz}^T] = \text{Covariance of fMRI responses}$$

of stimulus images



Links

Paper:
arXiv:2411.08197
sharvey@flatironinstitute.org



Authors:
@sarah-harvey.bsky.social
@lipshutz.bsky.social
@itsneuronal.bsky.social

1. Yamins, D., DiCarlo, J. Using goal-driven deep learning models to understand sensory cortex. *Nat Neurosci* 19, 356–365 (2016).
2. Raghu, M., et al. “Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability”. *NeurIPS* 30 (2017).
3. Williams, A. H., et al. “Generalized Shape Metrics on Neural Representations”. *NeurIPS* Vol. 34, (2021).
4. Kriegeskorte, N., et al. “Representational similarity analysis—connecting the branches of systems neuroscience”. *Front Syst Neurosci*. 2008; 2:4.
5. Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. “Similarity of Neural Network Representations Revisited”. *ICML*, Vol. 97, (2019).
6. Muzellec, B., and Cuturi, M. “Generalizing Point Embeddings using the Wasserstein Space of Elliptical Distributions”. *NeurIPS* (2018).