

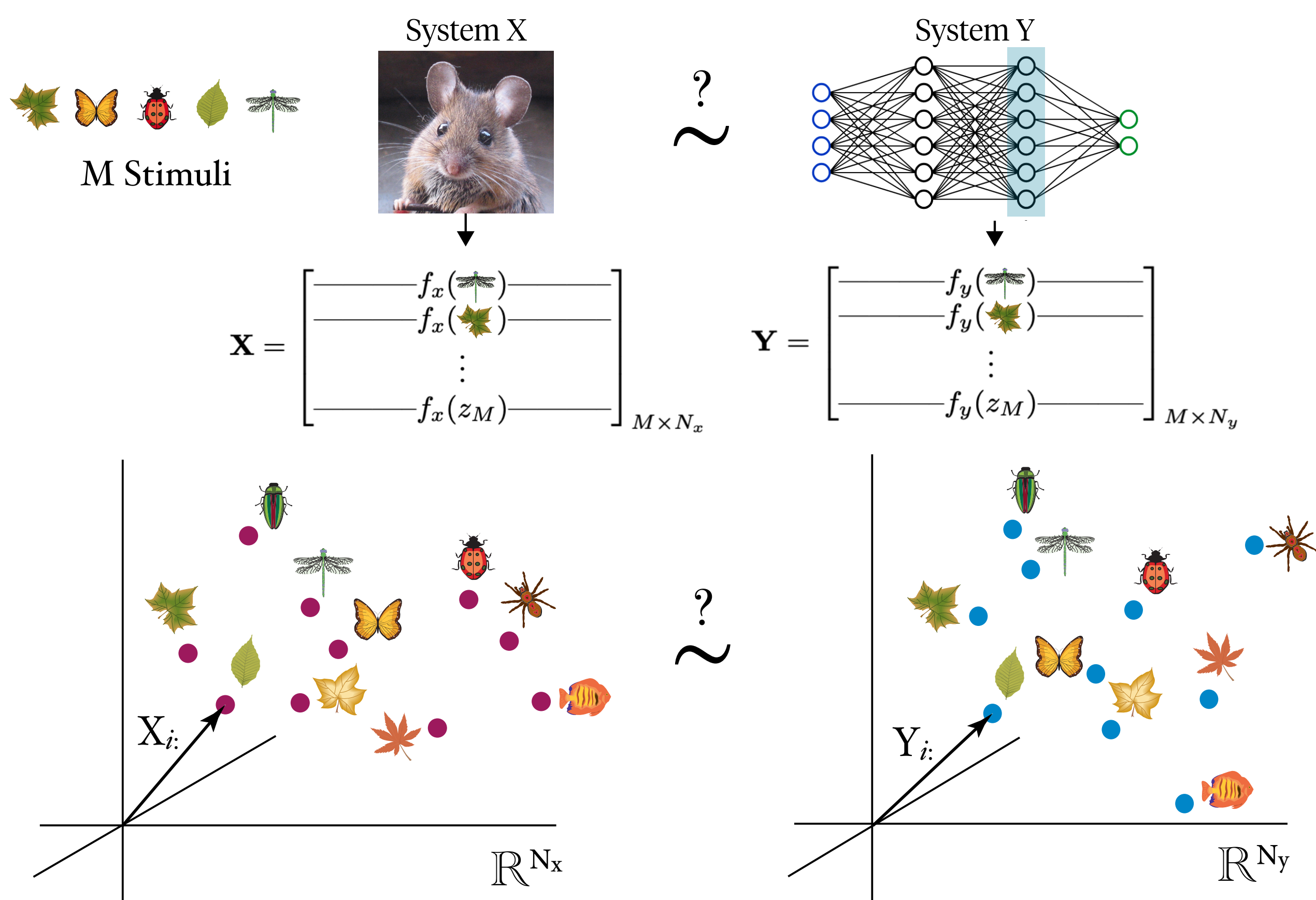
What Representational Similarity Measures Imply about Decodable Information

¹ Flatiron Institute, Center for Computational Neuroscience ² New York University, Center for Neural Science

Measuring Representational Similarity

Comparative analyses are important tools for understanding complex systems

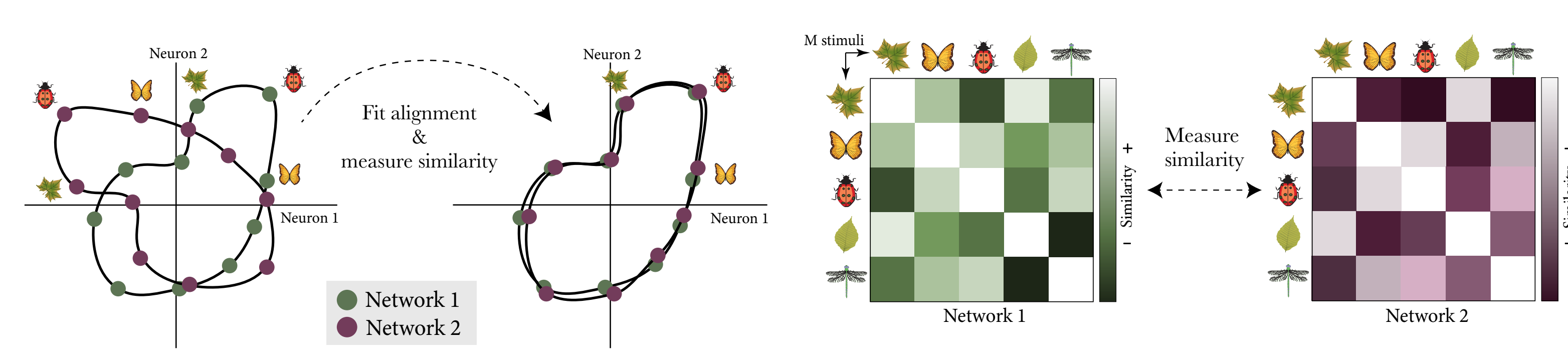
How do we quantify similarity between *neural representations*?



Many Existing Approaches

(Dis)similarity measures that transform or align neural dimensions

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



- Linear regression [1]
- Canonical correlations analysis [2]
- Procrustes shape distances [3]

- Representational Similarity Analysis (RSA) [4]
- Bures distance (Muzellec et. al. 2018) [6]
- Centered Kernel Alignment (CKA) [5]

$$\mathcal{P}(X, Y) = \min_{O \in \mathcal{O}(N)} \|\mathbf{X} - \mathbf{Y}O\|_F$$

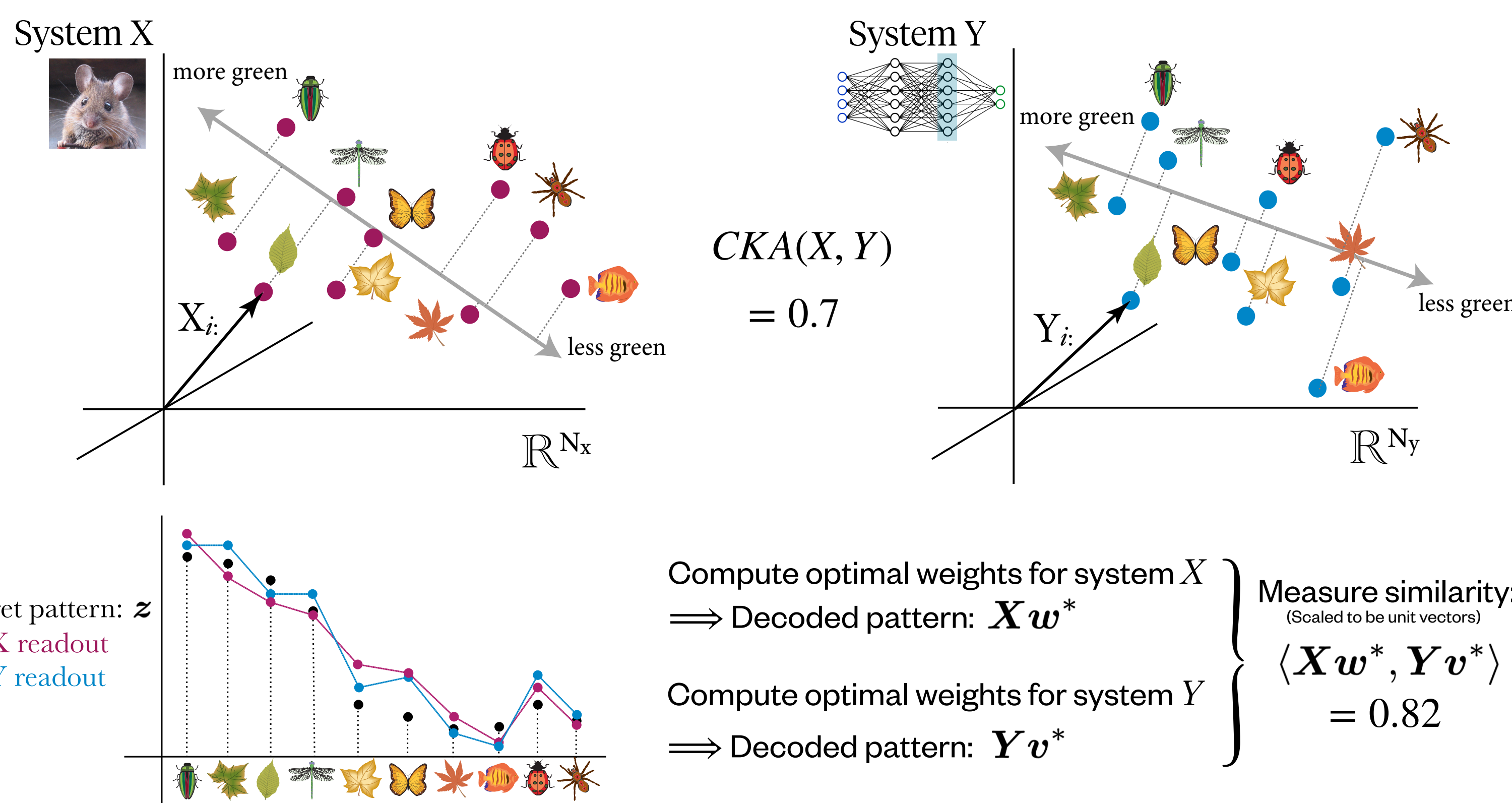
Measure Euclidean distance between representations after rotational alignment
Takes values from 0 to ∞ (distance)

$$CKA(X, Y) = \frac{\text{Tr} \mathbf{X}\mathbf{X}^T \mathbf{Y}\mathbf{Y}^T \text{ (Linear kernels)}}{\sqrt{\text{Tr}(\mathbf{X}\mathbf{X}^T)^2 \text{Tr}(\mathbf{Y}\mathbf{Y}^T)^2}}$$

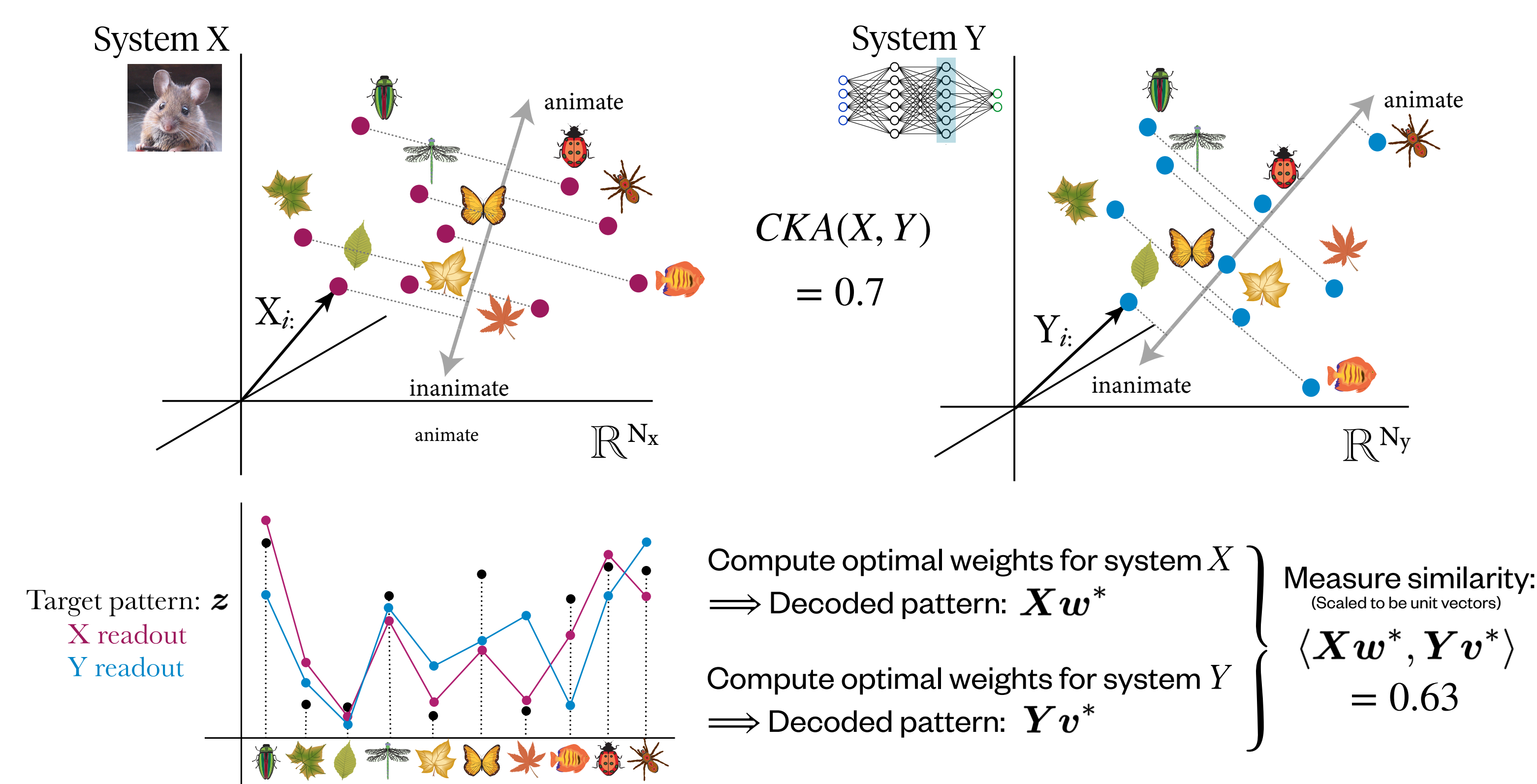
Compute inner product between centered representations and normalize
Takes values from 0 to 1 (similarity)

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

Idea: Measure similarity over an ensemble of decoding tasks

$$\text{average decoding similarity (ADS)} \quad \mathbb{E}_{z \sim P_z} \langle \mathbf{X}\mathbf{w}^*, \mathbf{Y}\mathbf{v}^* \rangle$$

$$\text{average decoding distance (ADD)} \quad \mathbb{E}_{z \sim P_z} \|\mathbf{X}\mathbf{w}^* - \mathbf{Y}\mathbf{v}^*\|_2^2$$

To compute these, we must choose:

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

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

General mathematical framework

Set up a family of linear decoding problems

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

This problem has a nice closed form solution: $\mathbf{w}^* = \frac{1}{M} \mathbf{G}(\mathbf{X})^{-1} \mathbf{X}^T \mathbf{z}$ *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

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

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

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

Take $a=0, b=1$ \Rightarrow $\frac{\mathbb{E} \langle \mathbf{X}\mathbf{w}^*, \mathbf{Y}\mathbf{v}^* \rangle}{\sqrt{\mathbb{E} \langle \mathbf{X}\mathbf{w}^*, \mathbf{X}\mathbf{w}^* \rangle \mathbb{E} \langle \mathbf{Y}\mathbf{v}^*, \mathbf{Y}\mathbf{v}^* \rangle}} = CKA(X, Y)$

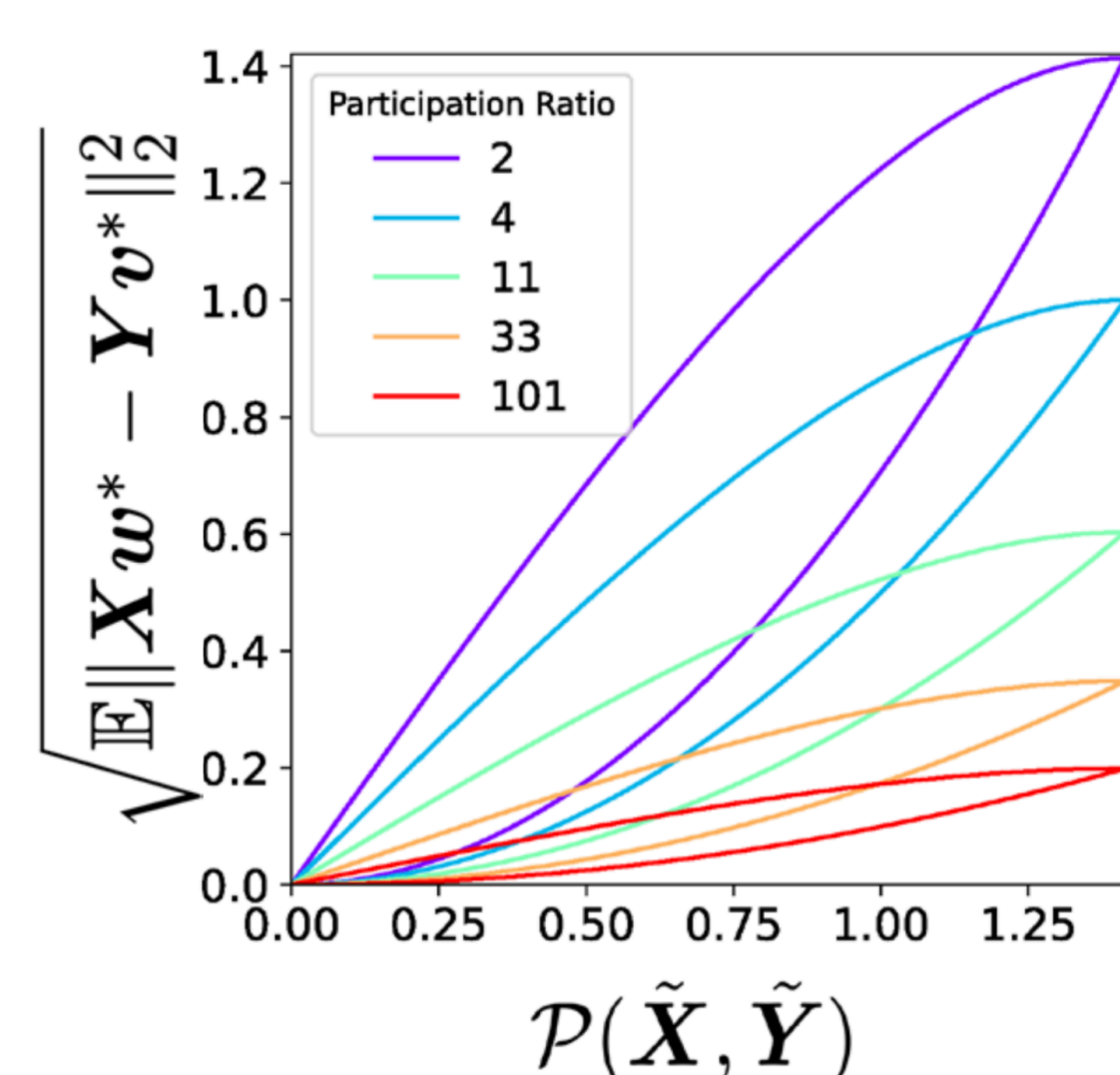
(Normalized) Average decoding similarity

More in paper:

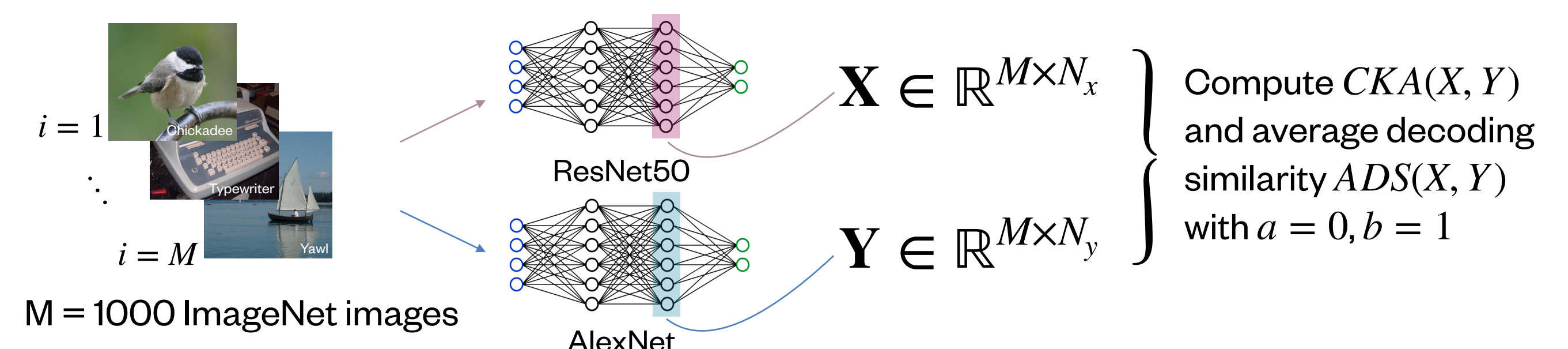
Similarity measure	a	b
Linear CKA	0	b
GULP	1	λ
CCA	1	0
ENSD	0	$\frac{1}{M} \frac{\text{Tr}[\mathbf{C}_X^2]}{\text{Tr}[\mathbf{C}_X]}$

We find saturated bounds relating average decoding distance to Procrustes shape distance between normalized representations

$\tilde{\mathbf{X}} := \frac{1}{\sqrt{M}} \mathbf{X}\mathbf{G}(\mathbf{X})^{-1/2}$ and $\tilde{\mathbf{Y}} := \frac{1}{\sqrt{M}} \mathbf{Y}\mathbf{G}(\mathbf{Y})^{-1/2}$

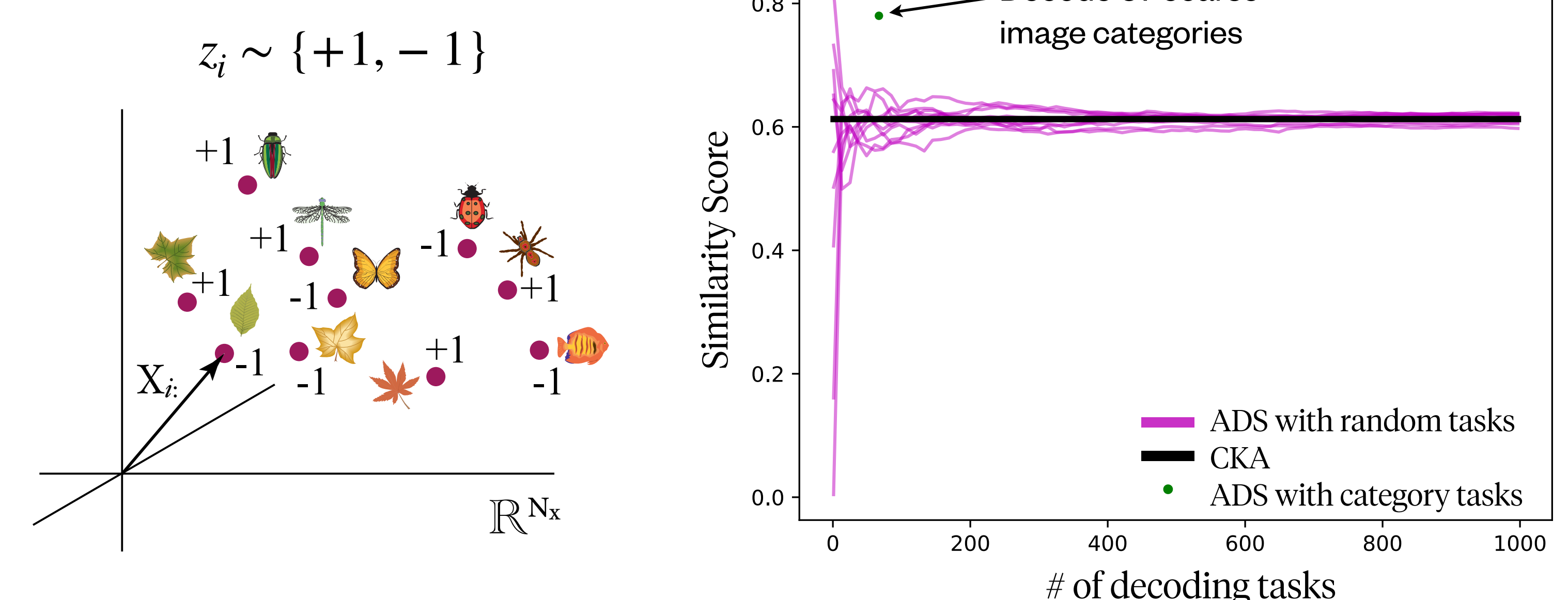


Example: CKA



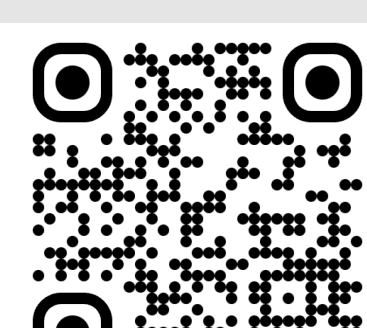
Ensemble of decoding tasks:

► Random binary partitions of input images



Links

Paper: arXiv:2411.08197
sharvey@flatironinstitute.org



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

References

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).