



**AI/CRV
2026**

NeurIPS 2025

**From Flat to Hierarchical: Extracting Sparse
Representations with Matching Pursuit**

**Valérie Costa^{1*} Thomas Fel^{2*} Ekdeep Singh Lubana^{3,4*}
Bahareh Tolooshams^{5,6†} Demba Ba^{2,7†}**

¹EPFL ²Kempner Institute, Harvard University

³CBS-NTT Program in Physics of Intelligence, Harvard University

⁴Physics of Artificial Intelligence Group, NTT Research, Inc., Sunnyvale, CA, USA

⁵University of Alberta ⁶Alberta Machine Intelligence Institute (Amii) ⁷SEAS, Harvard University

From Flat to Hierarchical: Extracting Sparse Representations with Matching Pursuit

Bahareh Tolooshams, Ph.D.

Assistant Professor,

University of Alberta

btolooshams@ualberta.ca

<https://btolooshams.github.io>

May 28, 2026



NeuBahar Lab
**UNIVERSITY
OF ALBERTA**

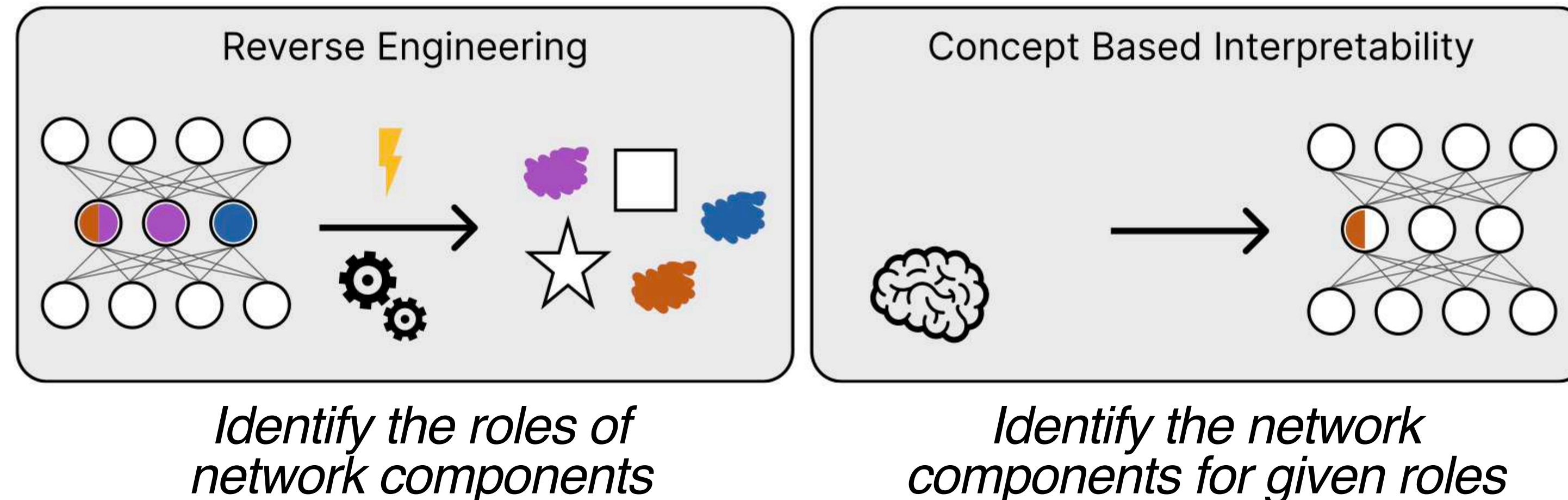


amii

CIFAR

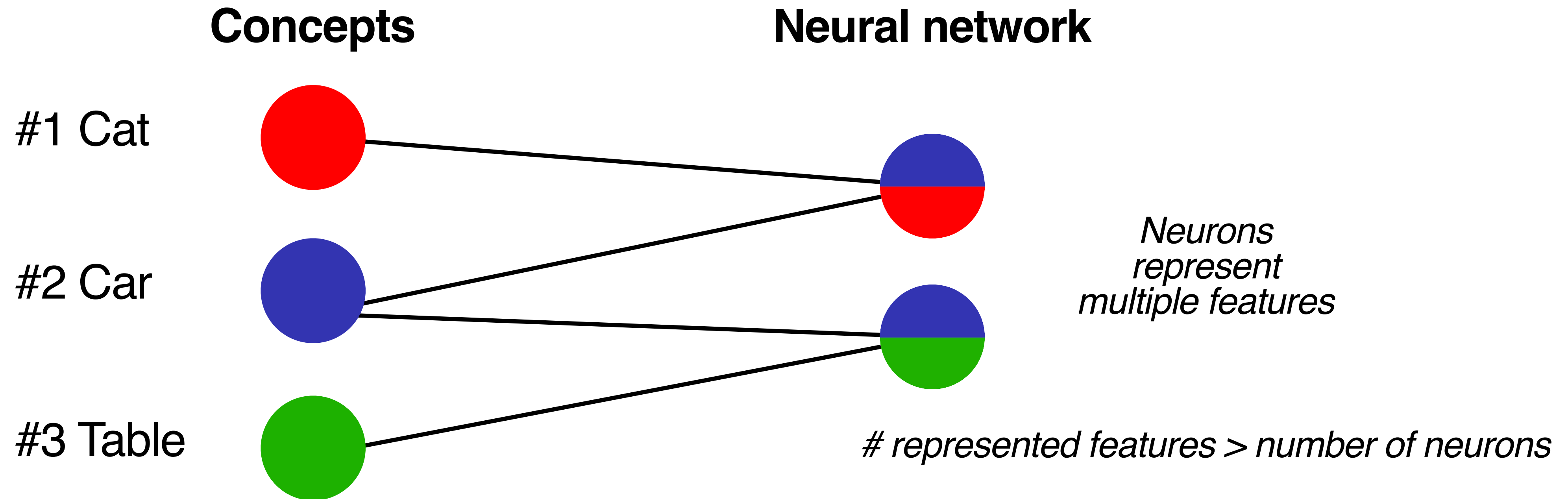
How Does Large Neural Networks Represent Data?

Mechanistic Interpretability: Understanding the structures underlying internal model representations.



- Decompose the internal AI representation into *interpretable* components.

How Does AI Represent Data?



- Features are sufficiently **sparse**, so assign features rarely co-occur to one neuron.
- As sparsity increases, models use **superposition** to represent more features than dimensions.

Linear Representation Hypothesis (Model of Data)

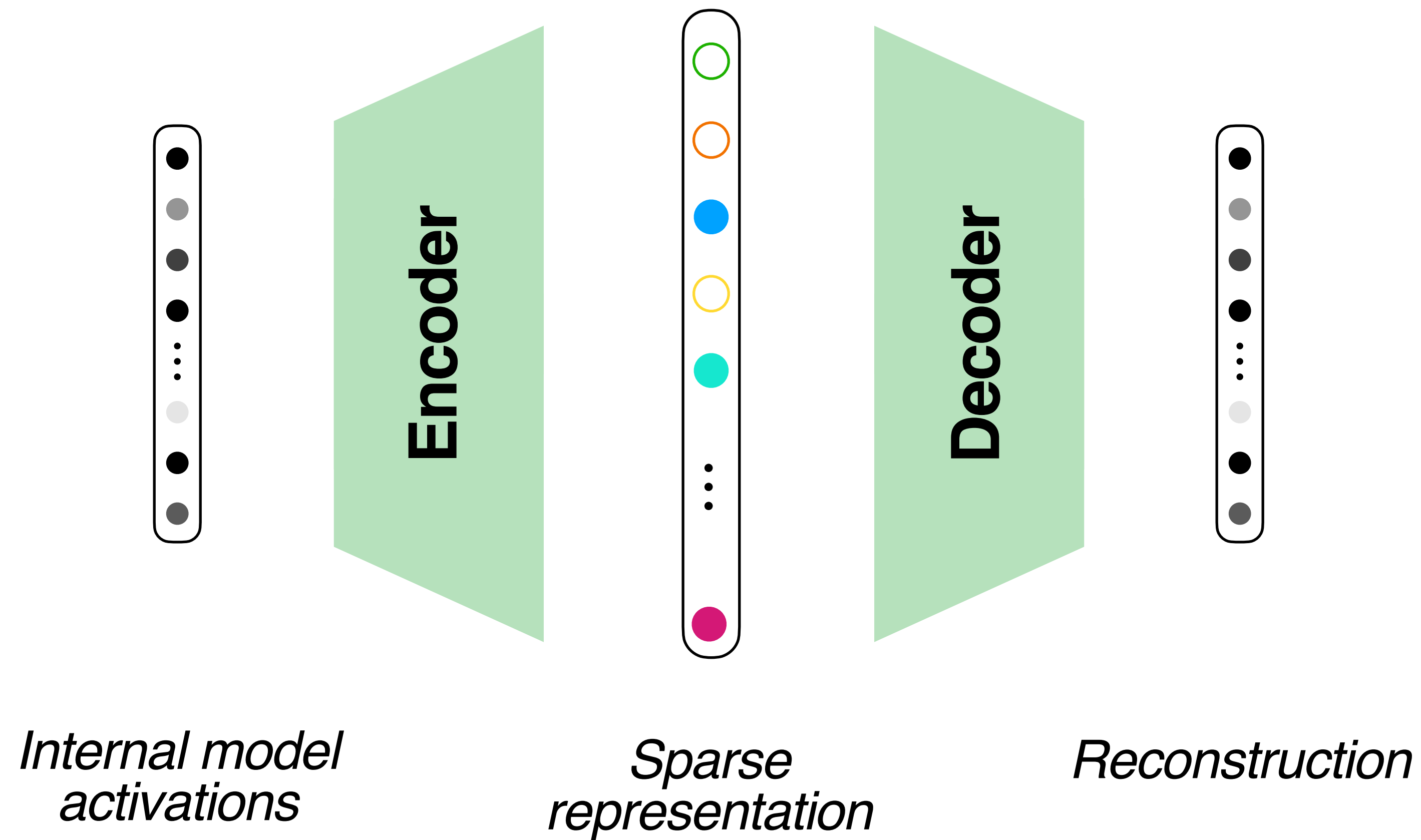
Definition 2.1 (Linear Representation Hypothesis (LRH)). A representation $\mathbf{x} \in \mathbb{R}^m$ is said to satisfy the linear representation hypothesis (LRH) if there exists a dictionary $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_p] \in \mathbb{R}^{m \times p}$ and a coefficient vector $\mathbf{z} \in \mathbb{R}^p$ such that $\mathbf{x} = \mathbf{D}\mathbf{z}$, under the following conditions:

$$\left\{ \begin{array}{ll} \text{(i) Overcompleteness:} & p \gg m; \\ \text{(ii) Quasi-orthogonality:} & \max_{i \neq j} |\mathbf{D}_i^\top \mathbf{D}_j| \leq \varepsilon, \quad \text{where } \forall i \quad \|\mathbf{D}_i\|_2 = 1; \text{ and} \\ \text{(iii) Sparsity:} & |\text{supp}(\mathbf{z})| \leq k \ll p. \end{array} \right.$$

- Features are represented by **directions**.
- Features are linearly **accessible**.

We have “model of data”; we can now construct a sparse autoencoder (SAEs) to extract geometry of representations in large language models.

Sparse Autoencoders (SAEs) to Extract Concepts from Model Representations



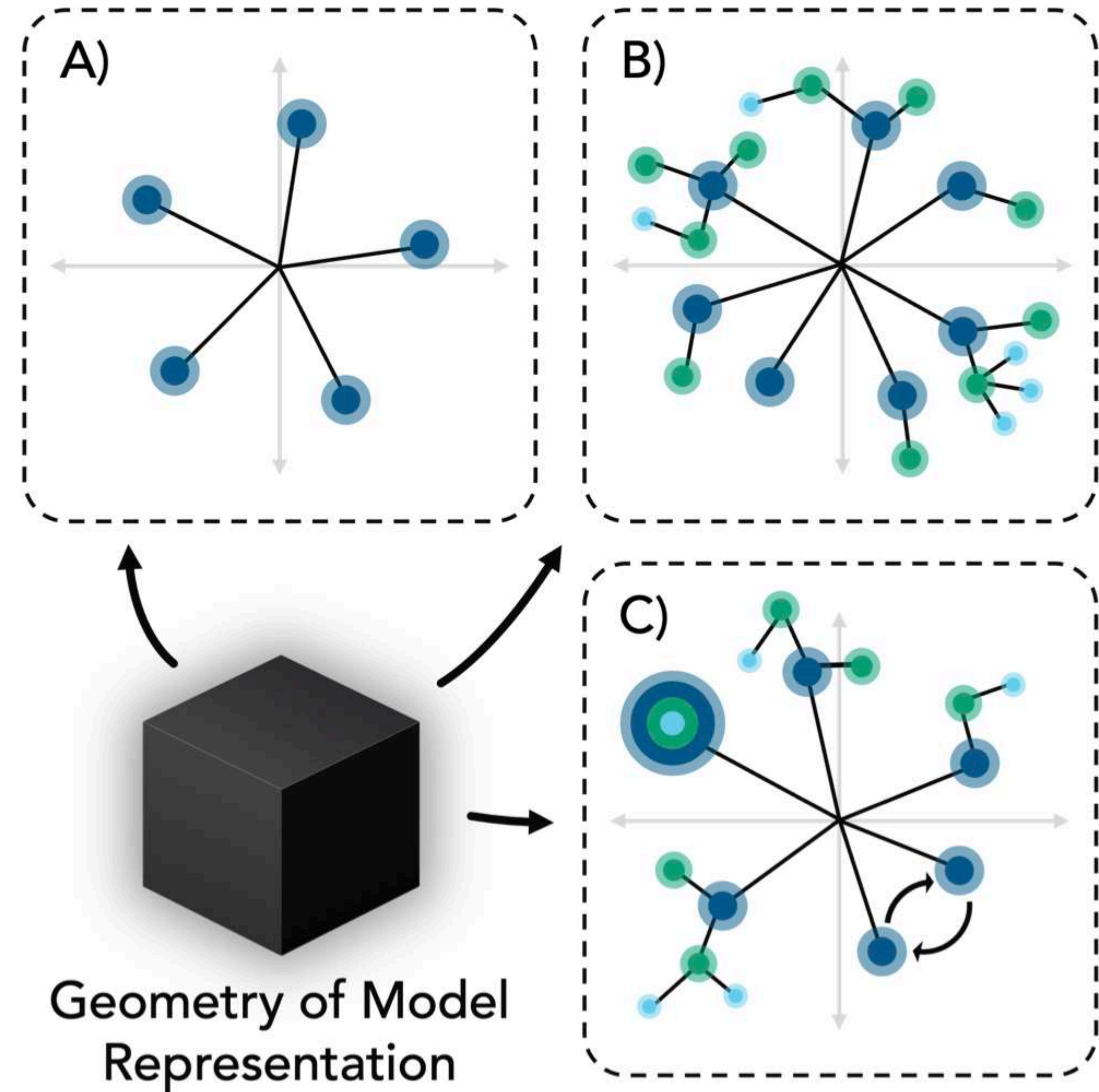
Can SAEs recover features beyond the *Linear Representation Hypothesis*?

Conceptual Organization in Neural Representations

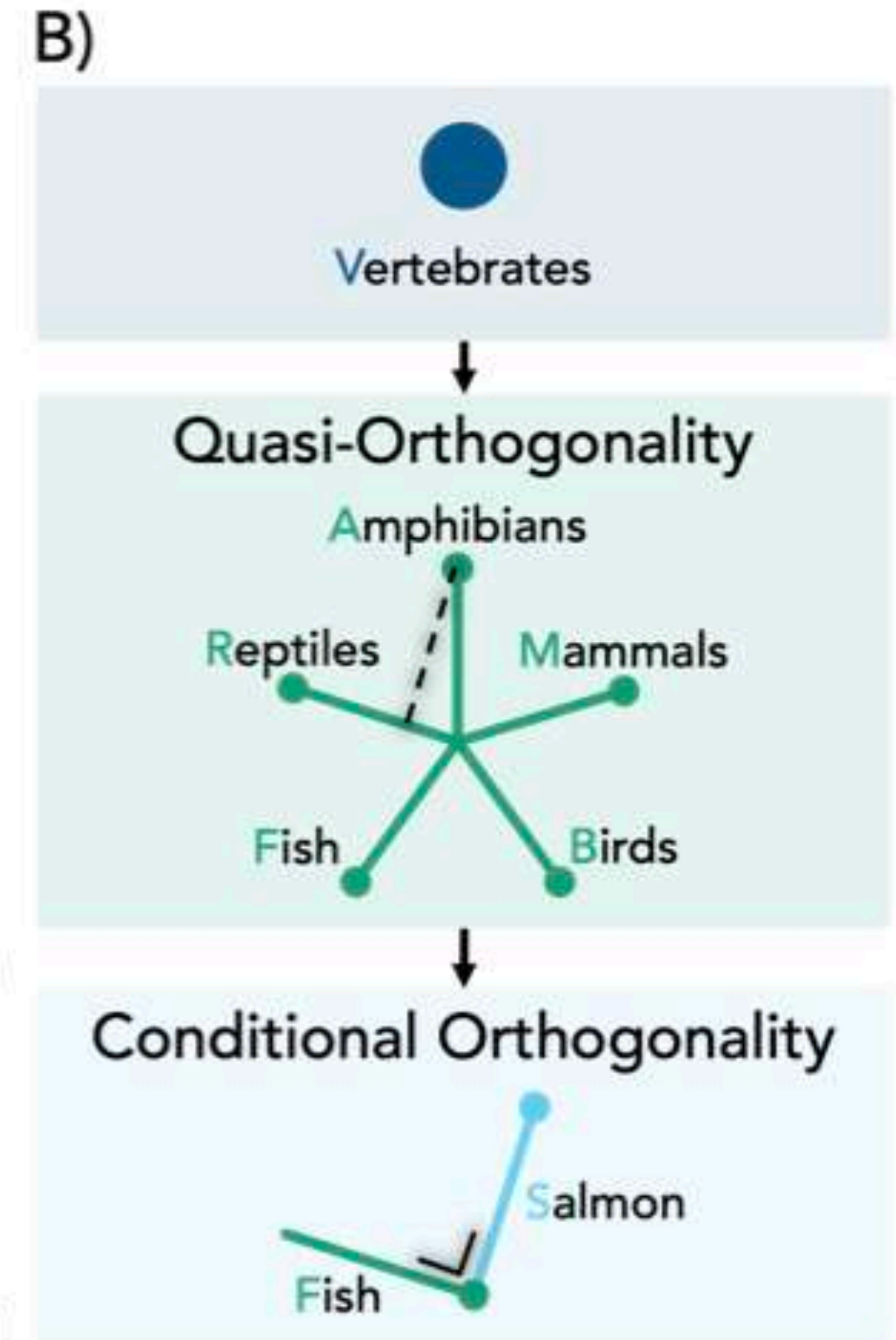
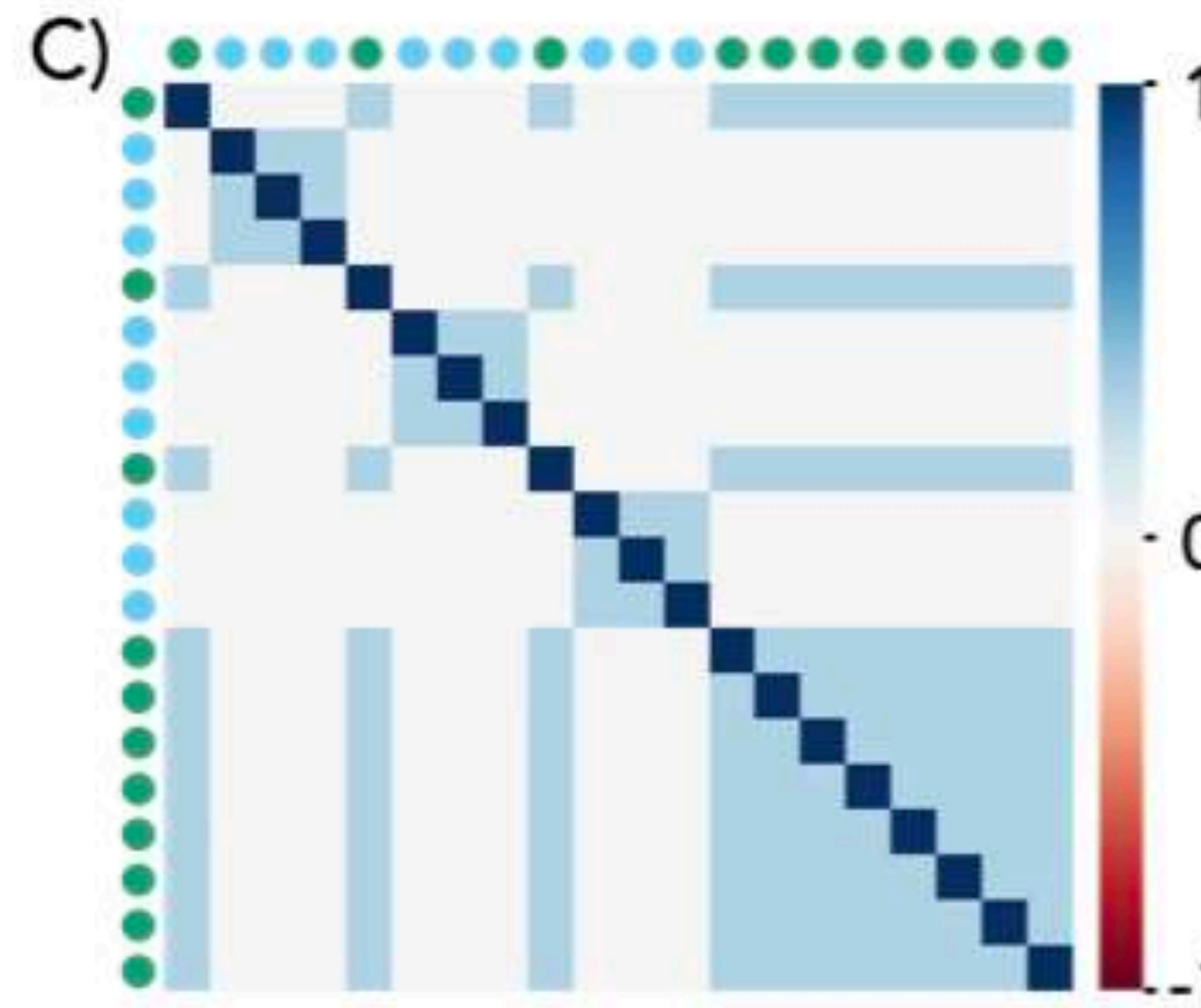
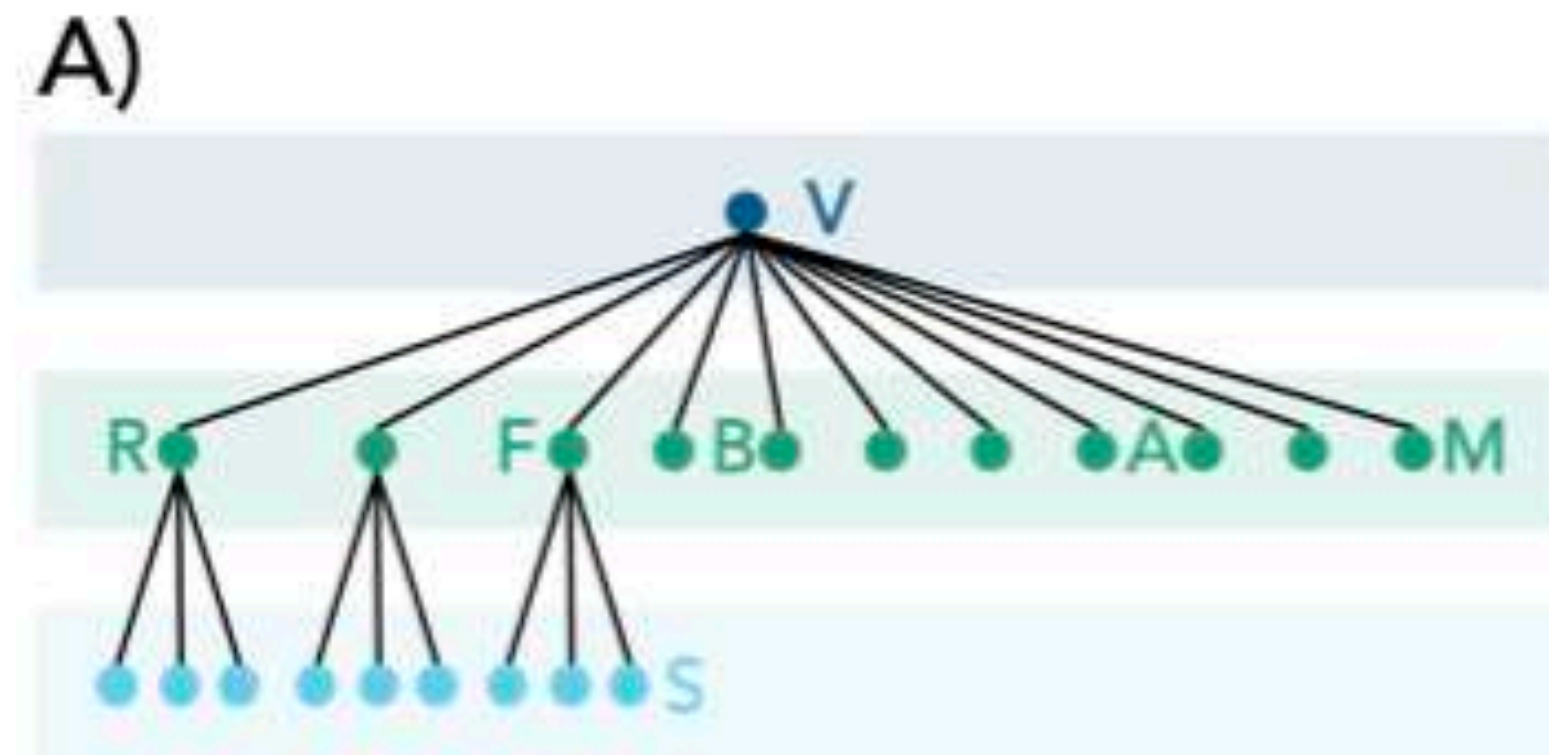
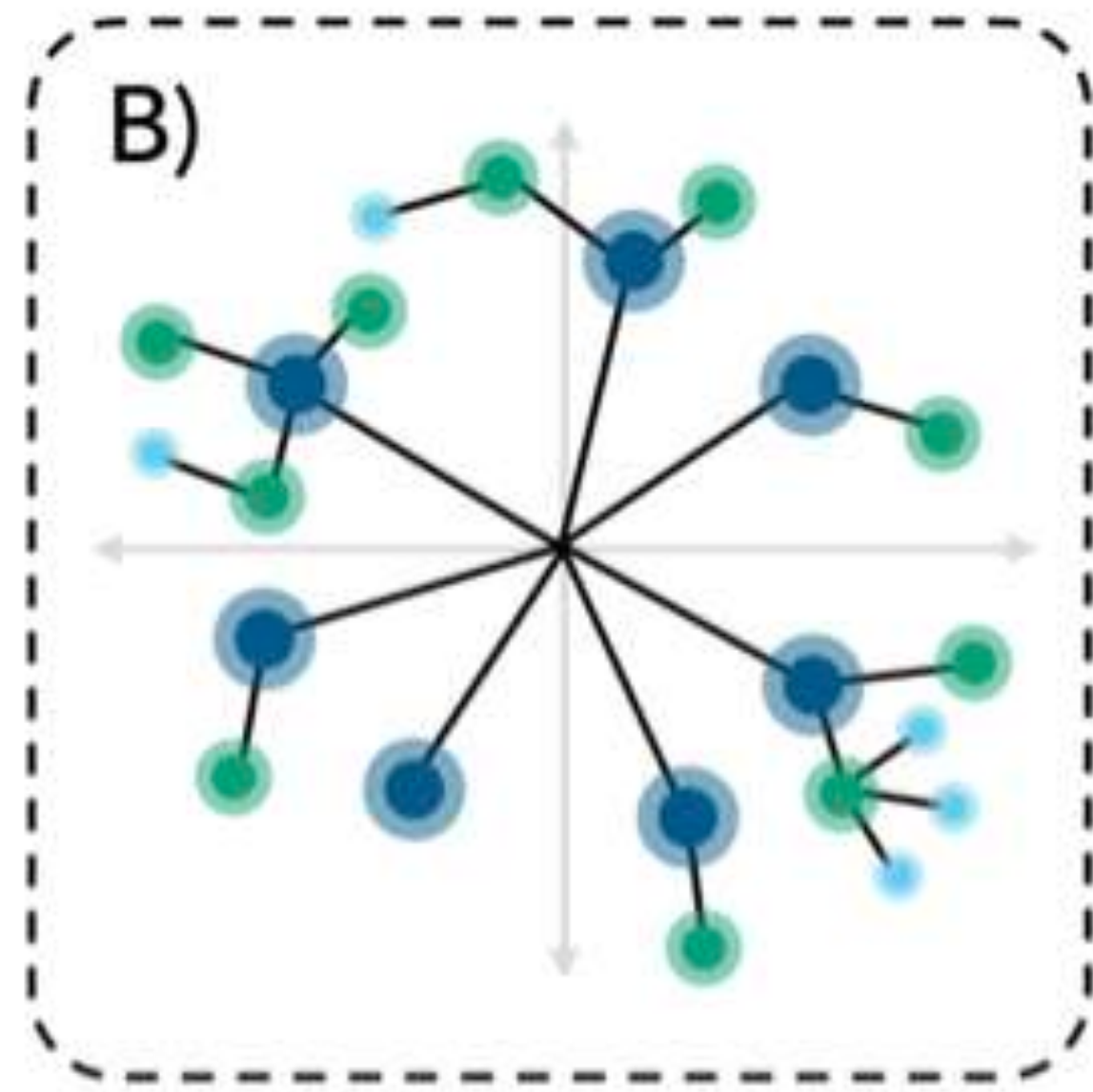
A) **Linearly accessible**: approximately orthogonal abstract directions, a.k.a Linear Representation Hypothesis (LRH).

B) **Hierarchical**: representations are structured in parent-child relations.

C) **Nonlinear, multidimensional, and temporally structured ...**



From Flat to Hierarchical Representations

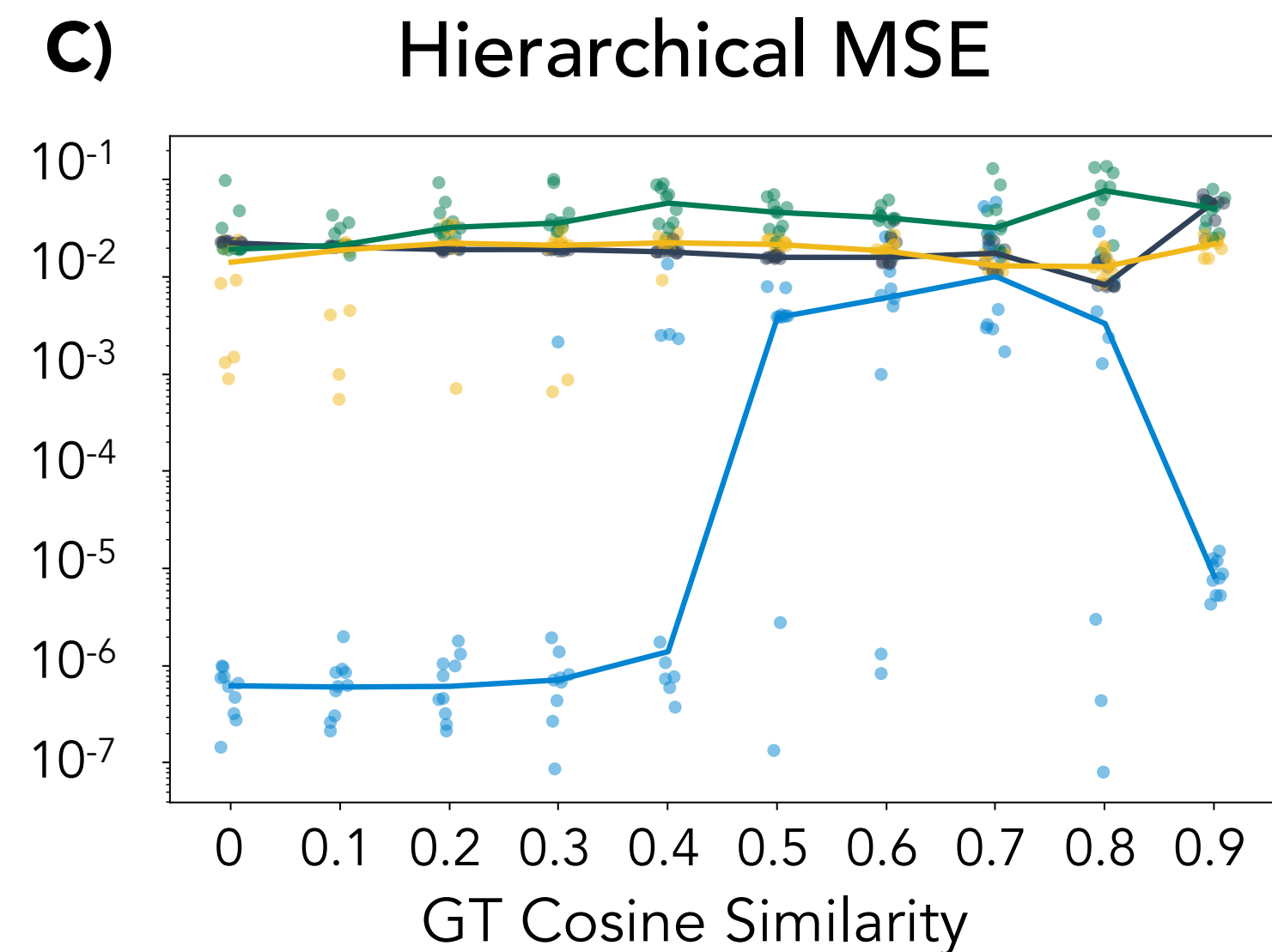
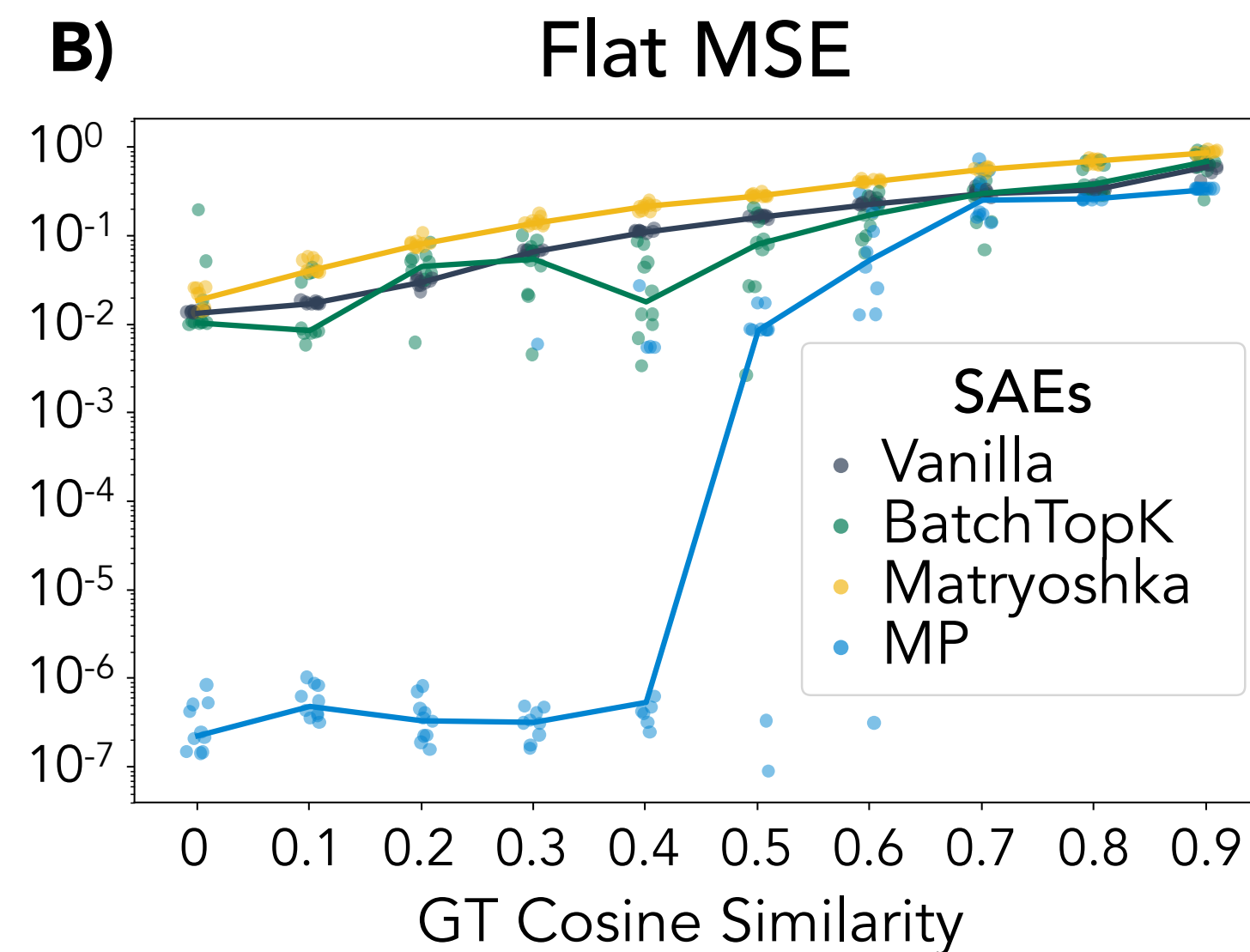
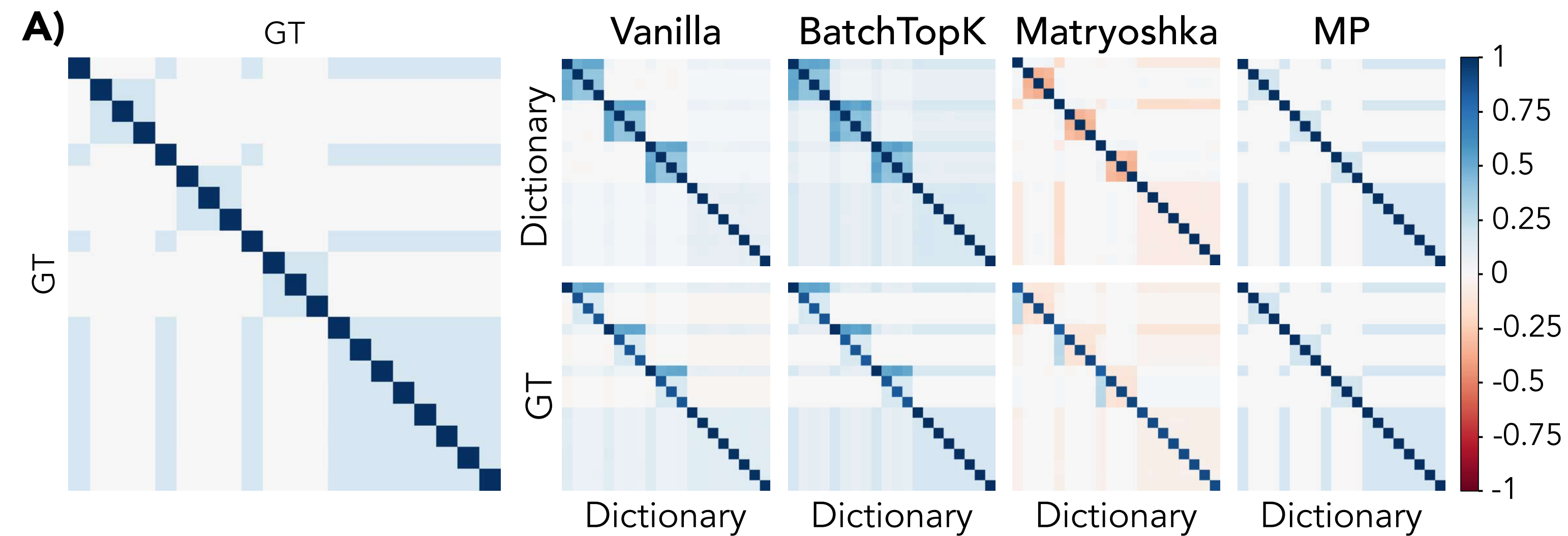


Extracting Hierarchical Representations with Matching Pursuit Sparse Autoencoder (MP-SAE)



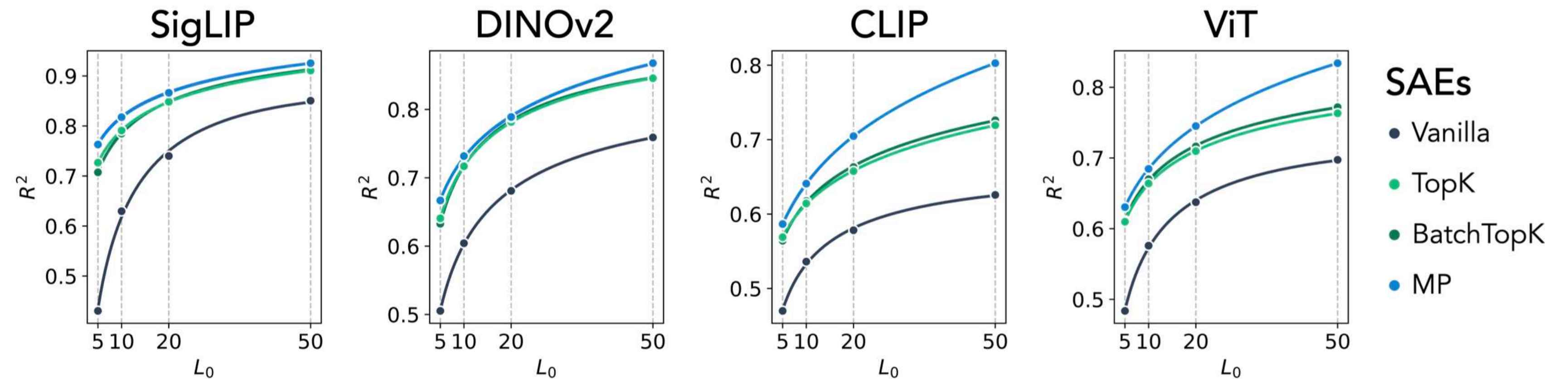
$$\mathbf{x} = \mathbf{r}^{(0)} = \underbrace{\varphi(\mathbf{x})}_{\text{linearly accessible}} + \underbrace{\sum_{t=1}^T \varphi(\mathbf{r}^{(t)})}_{\text{nonlinearly accessible}} + \underbrace{\mathbf{r}^{(T+1)}}_{\text{residual error}},$$

MP-SAE Recovers Hierarchical Correlated Features

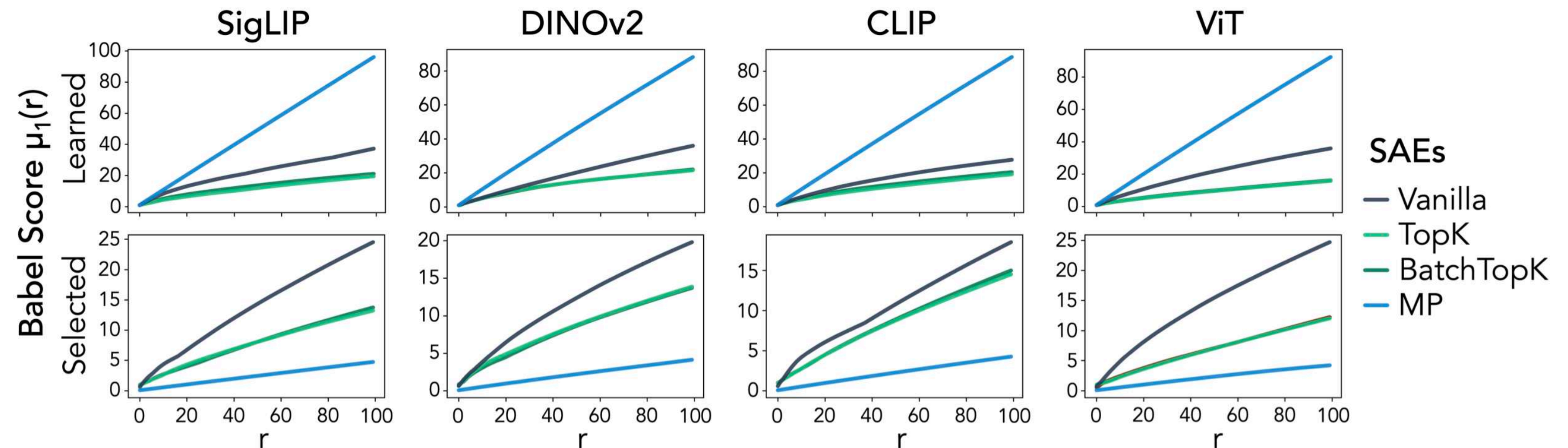


Extracting Features from Internal Representations of Vision Models

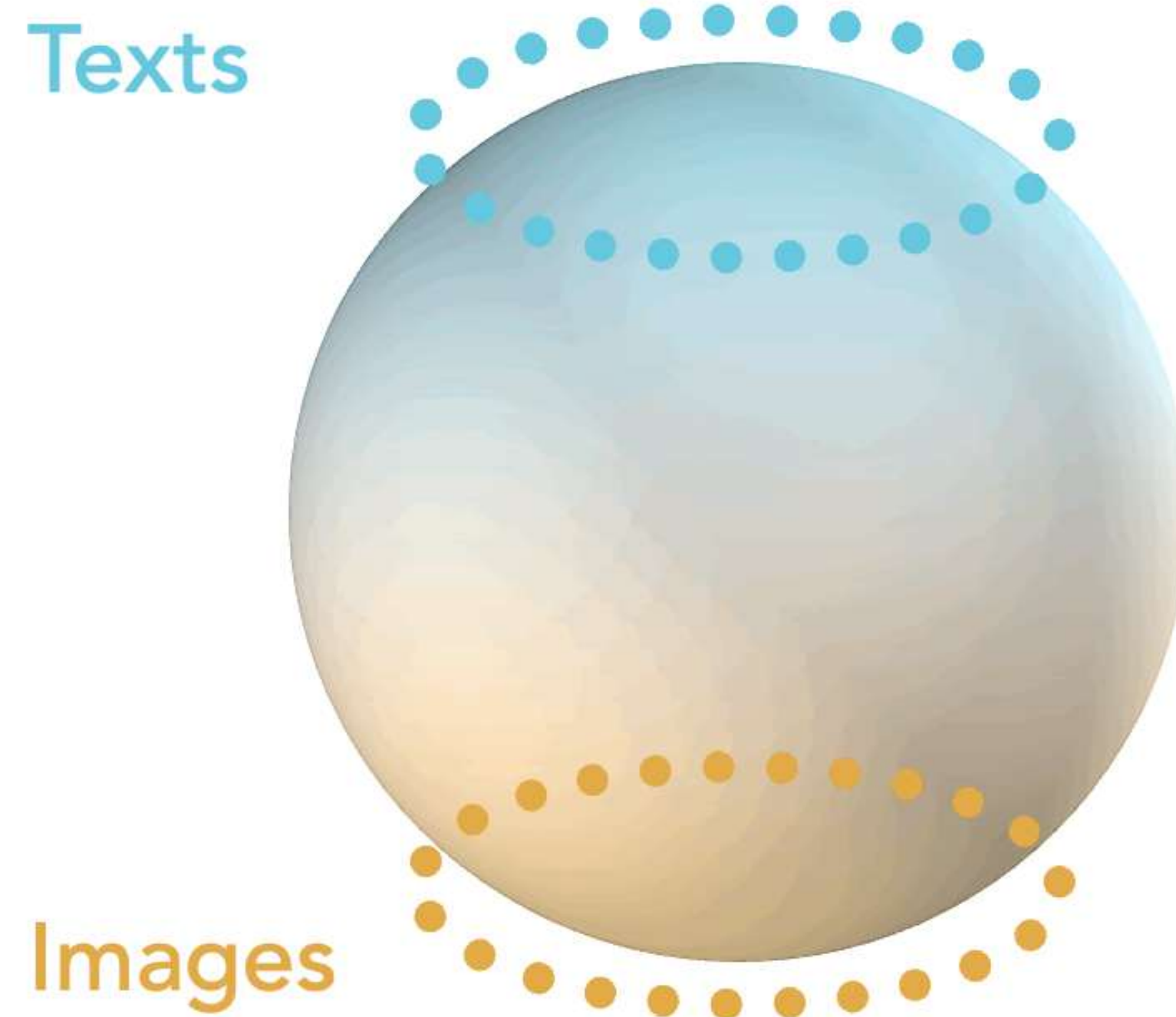
MP-SAE recovers more expressive features than standard SAEs.



MP-SAE promotes conditional orthogonality at inference.

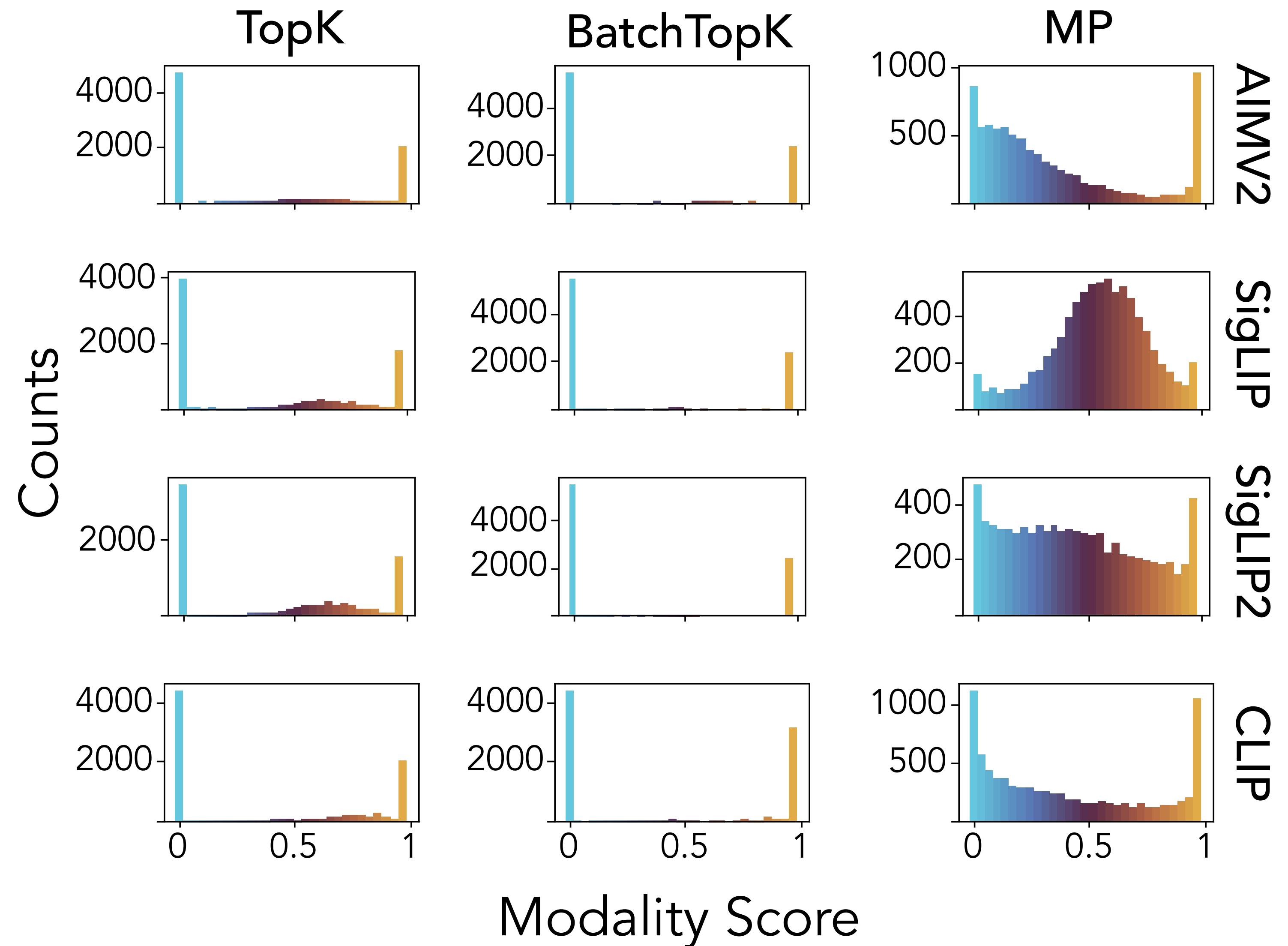


Interpreting Multi-Modal Large Models

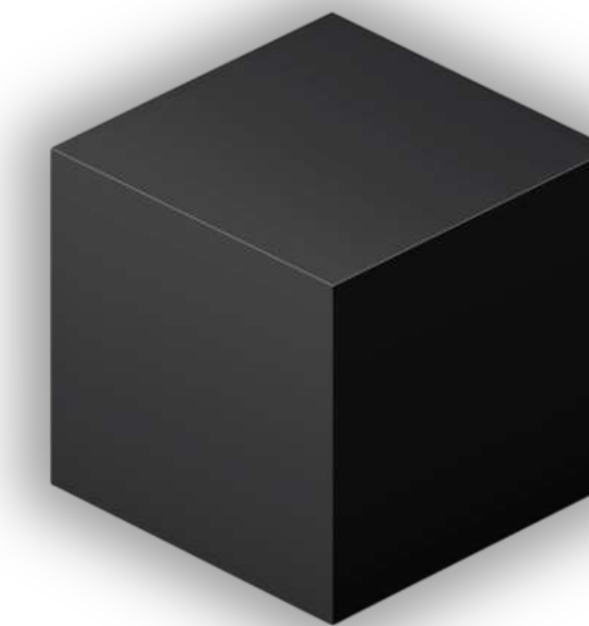
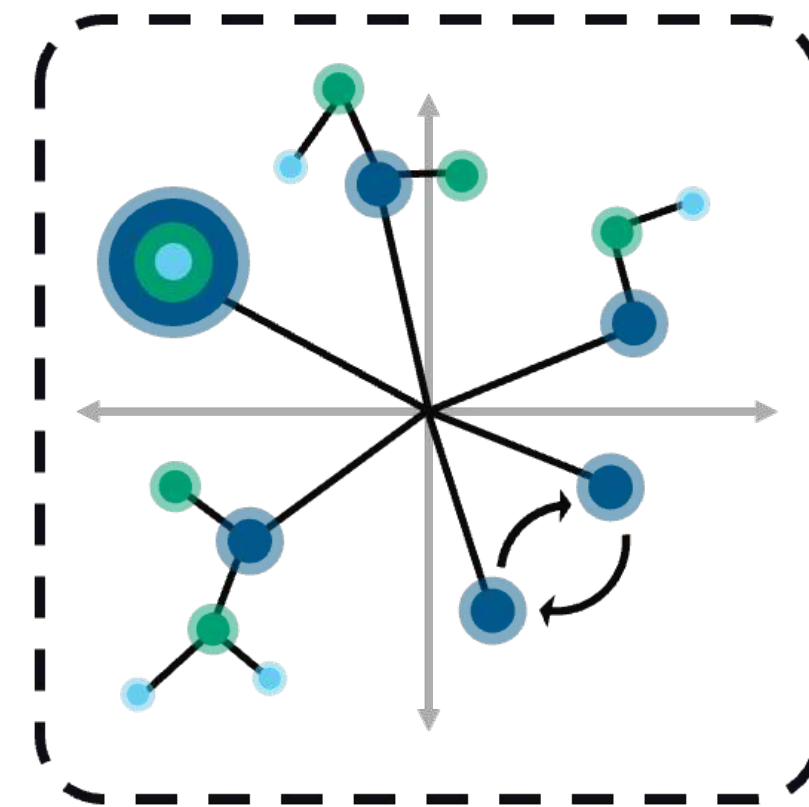
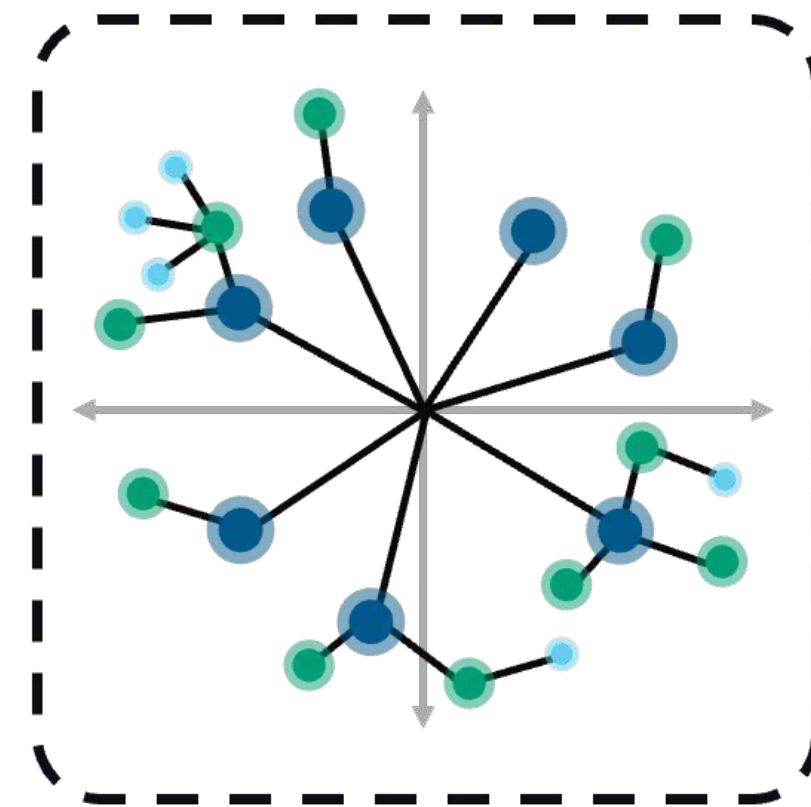
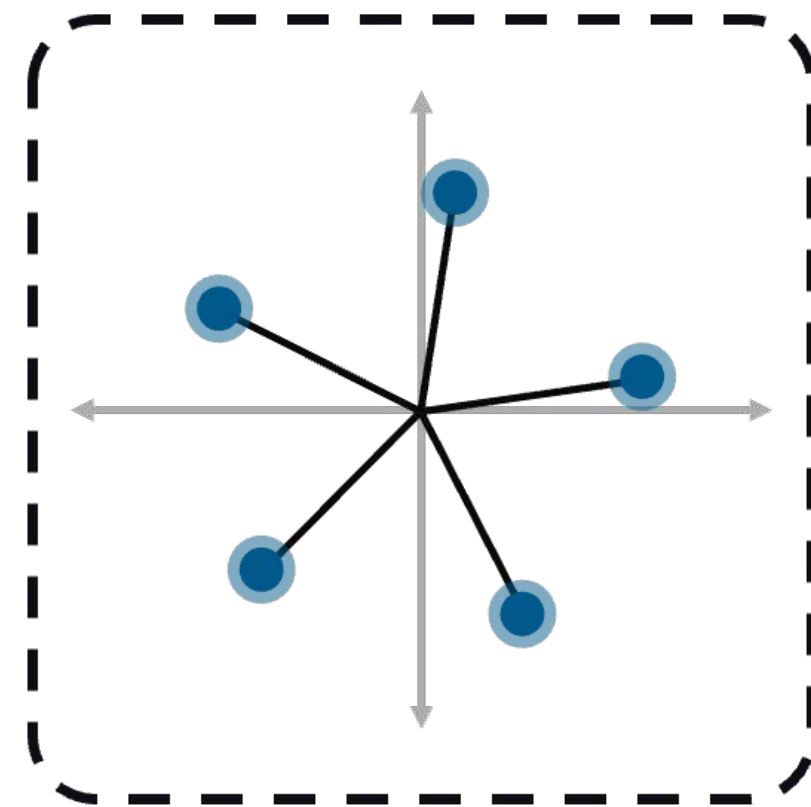


Identifying Shared Structures Across Modalities

Emergent modality subtraction from abstract meaning



Interpretability should begin with **geometry** of model **representations**, with methods to **reflect** those **assumptions**.



Geometry
of Model
Representation ?

Valérie Costa



Thomas Fel



Ekdeep Singh Lubana



Demba Ba



HARVARD
UNIVERSITY



Kempner
INSTITUTE

EPFL



NTTResearch



NeuBahar Lab
UNIVERSITY
OF ALBERTA



ami

Bertarelli Foundation

CIFAR