

From Flat to Hierarchical: Extracting Sparse Representations with Matching Pursuit

NeurIPS 2025

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

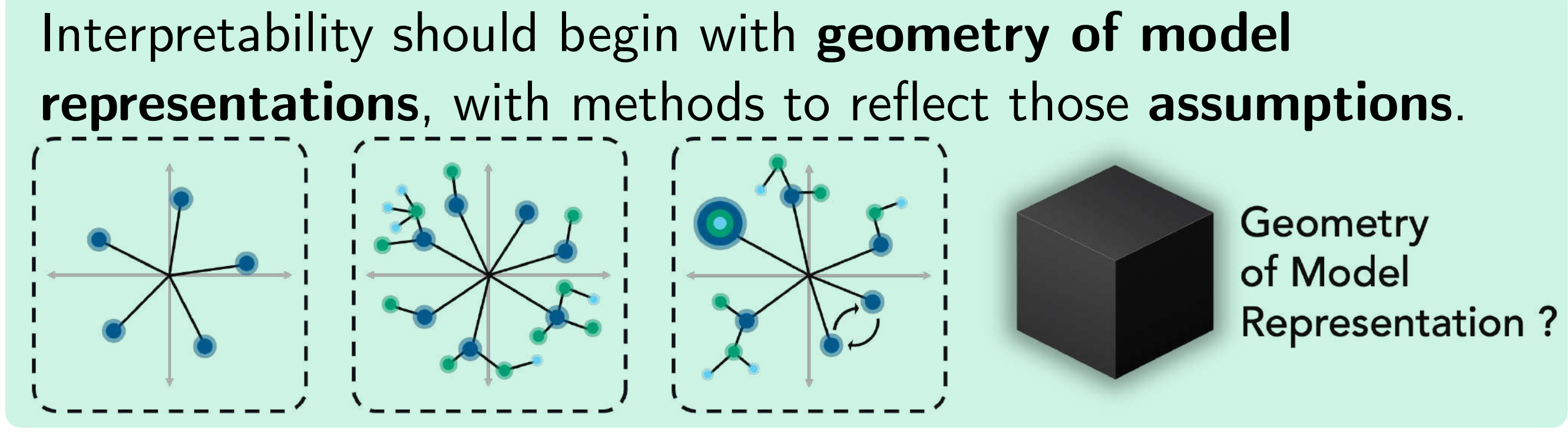
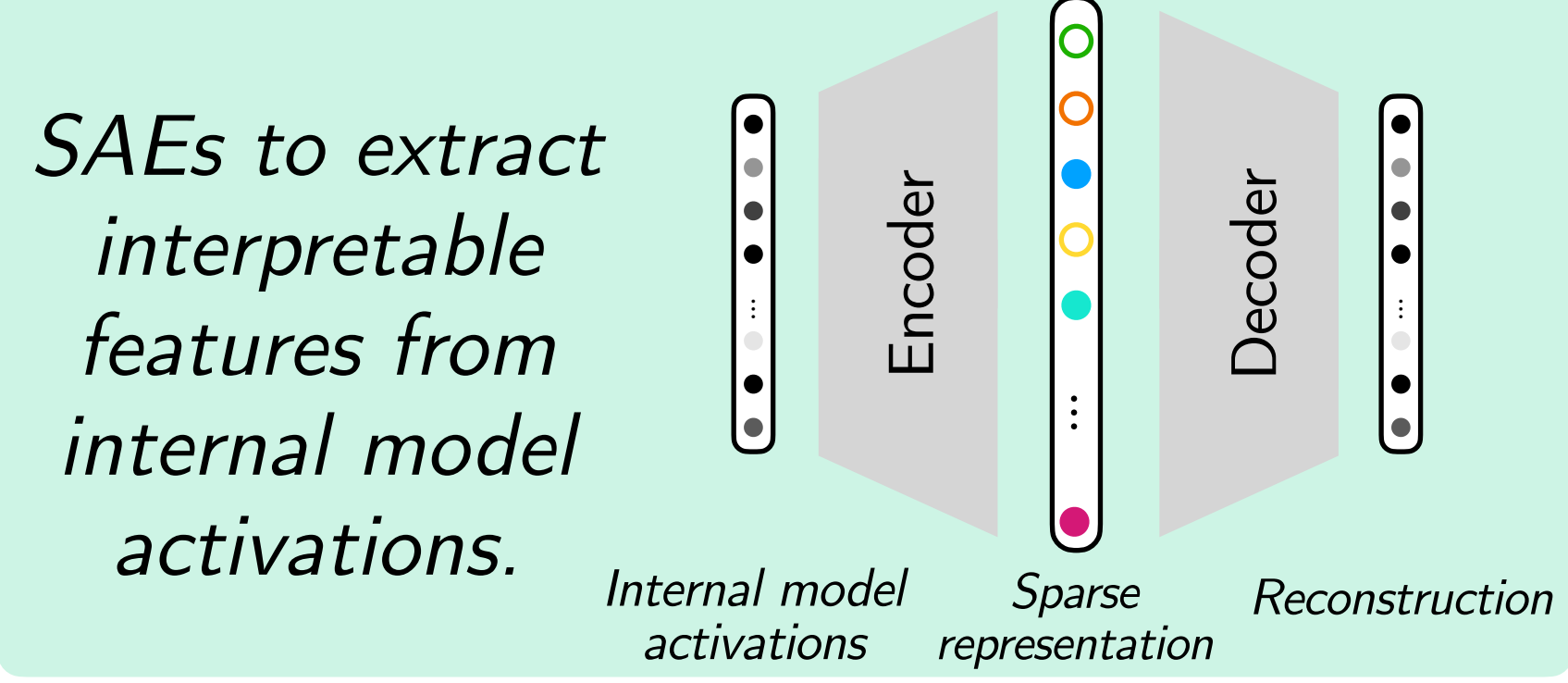
Nectar spotlight AI/CRV 2026



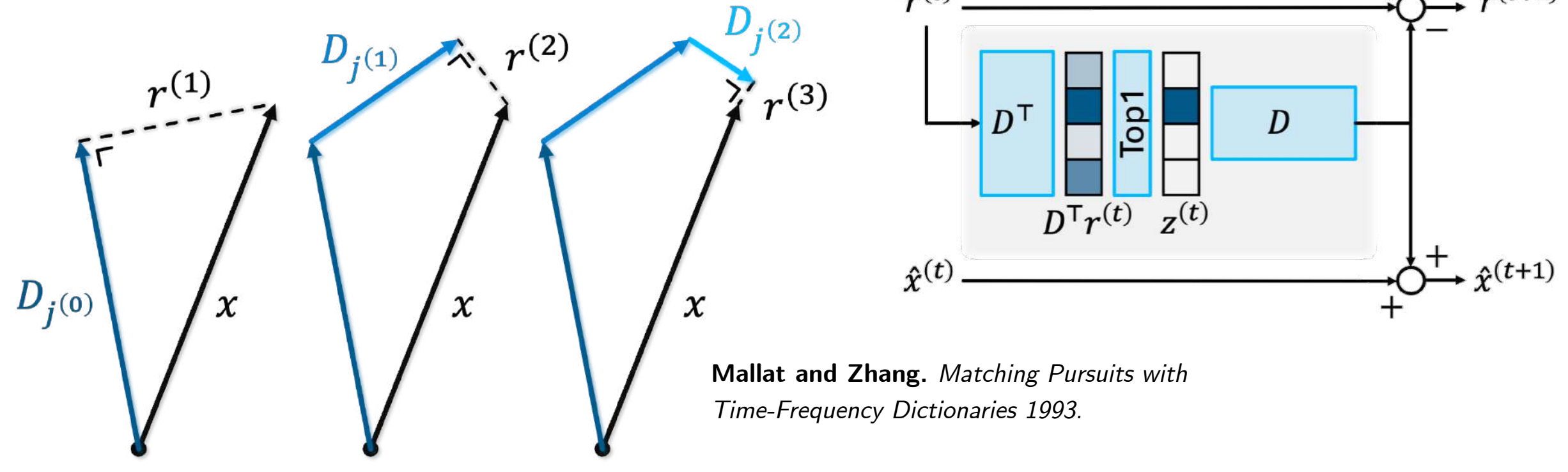
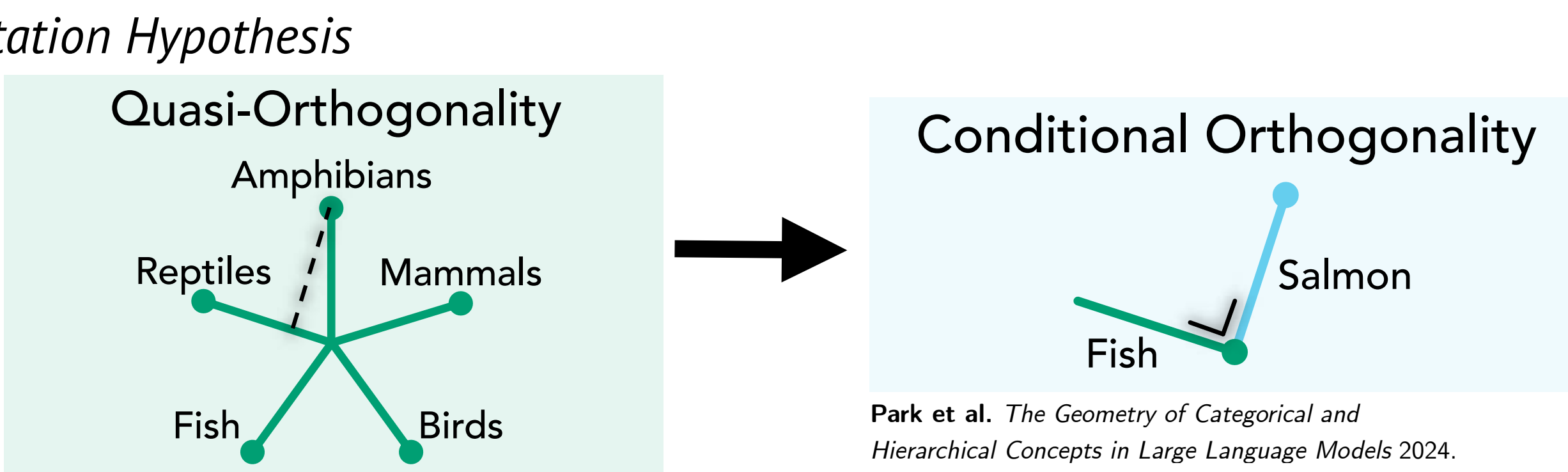
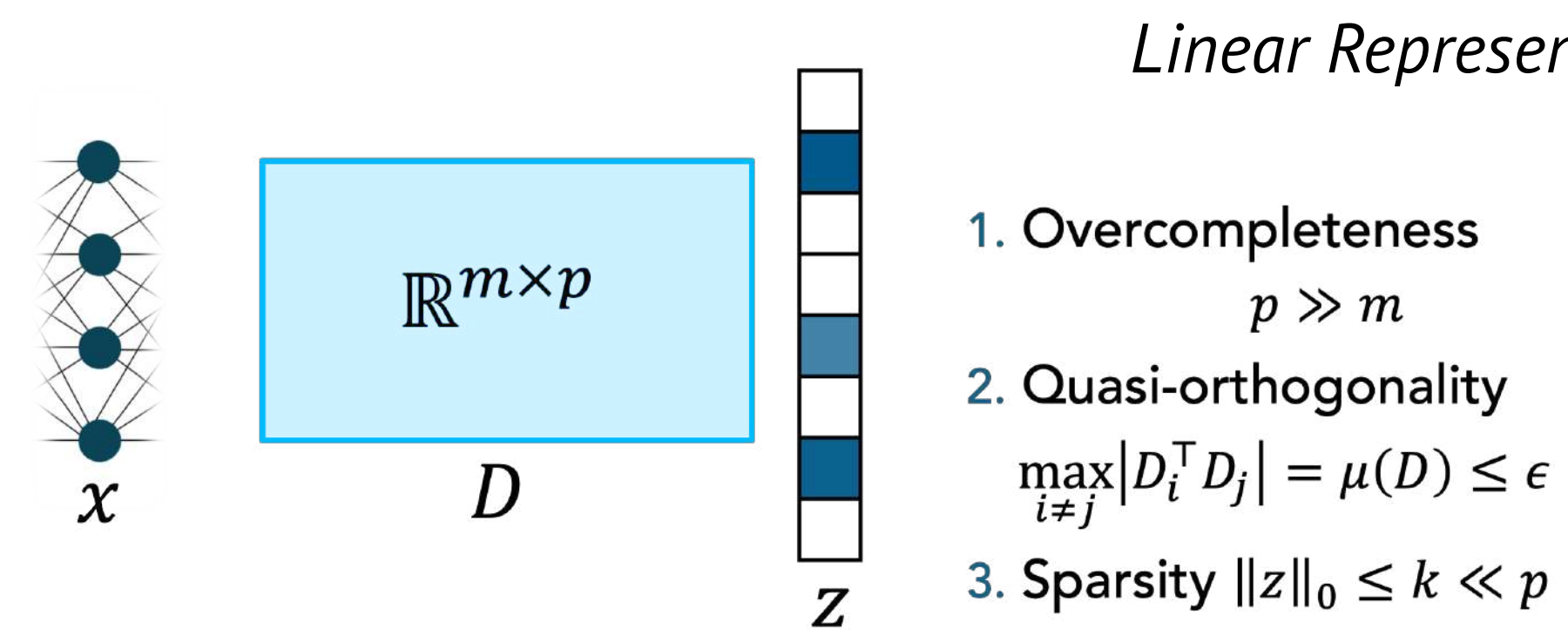
Motivated by the linear representation hypothesis (LRH), **sparse autoencoders (SAEs)** have become a popular tool for **mechanistic interpretability**.

Can SAEs recover features beyond the LRH?

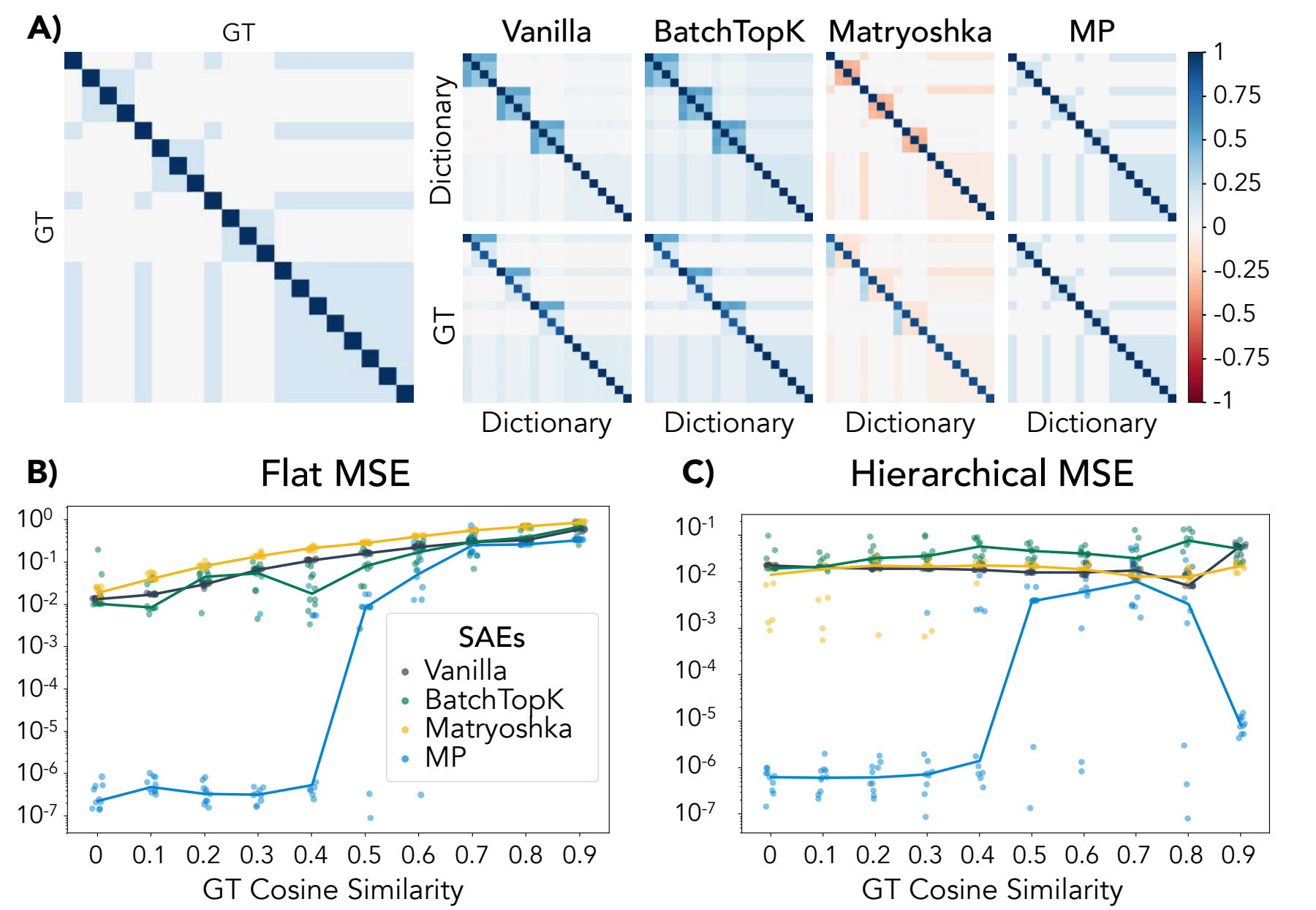
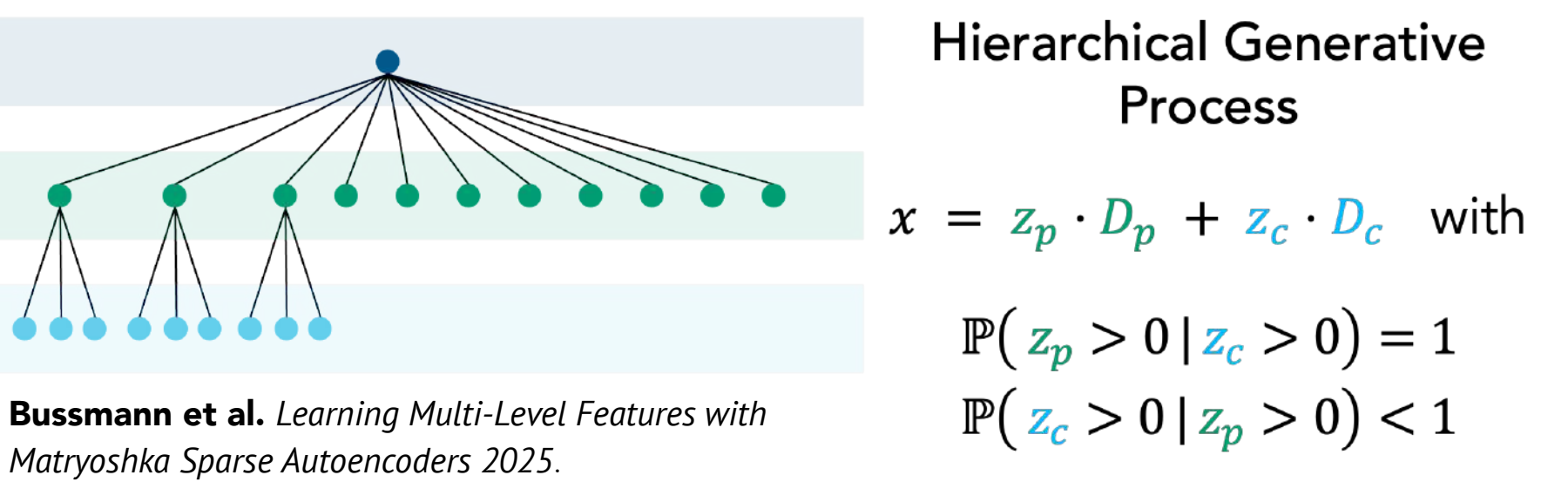
We propose **MP-SAE** to extract hierarchical features.



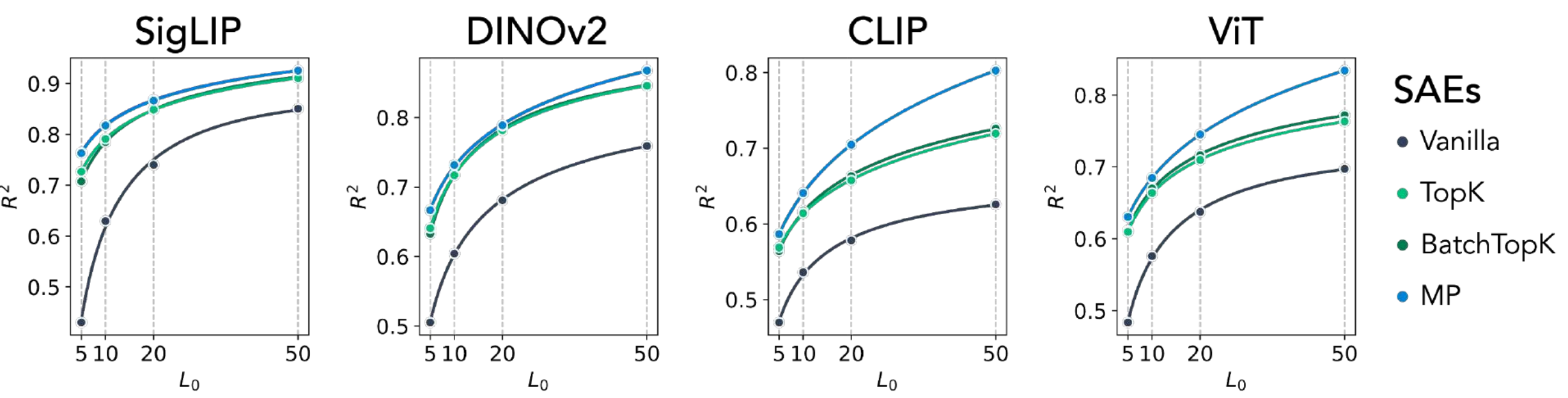
Method: Conditional Orthogonality and Matching Pursuit Sparse Autoencoder (MP-SAE)



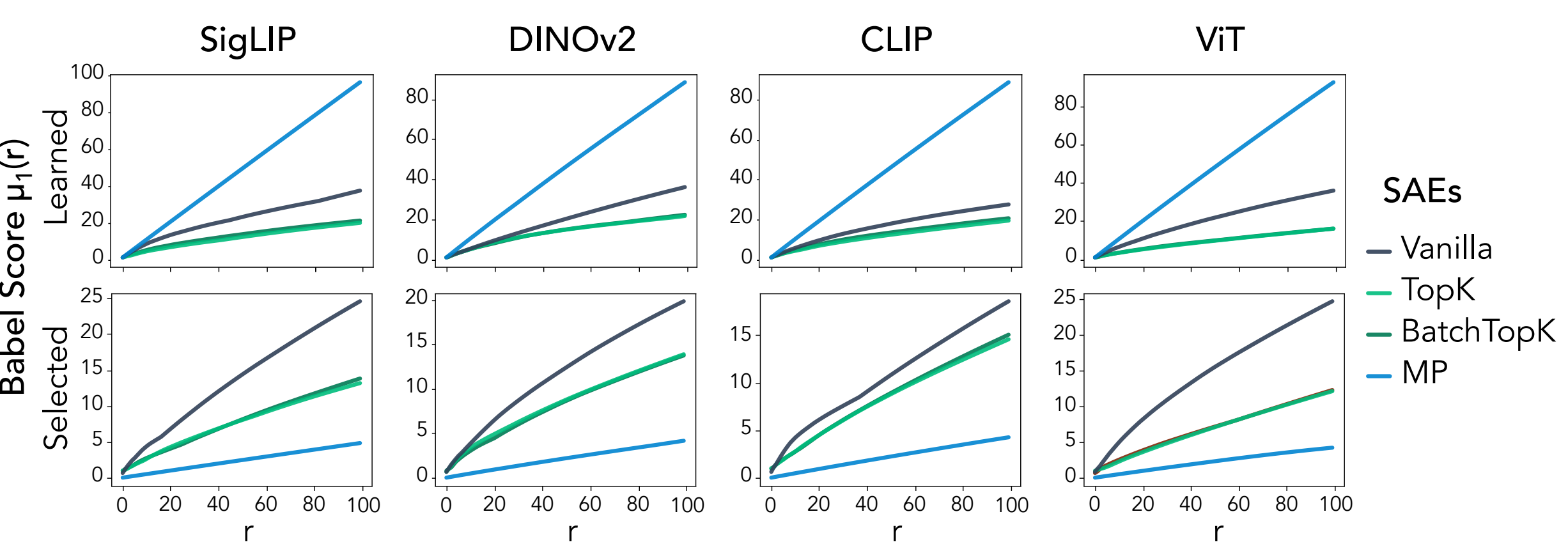
MP-SAE Recovers Hierarchical Correlated Features



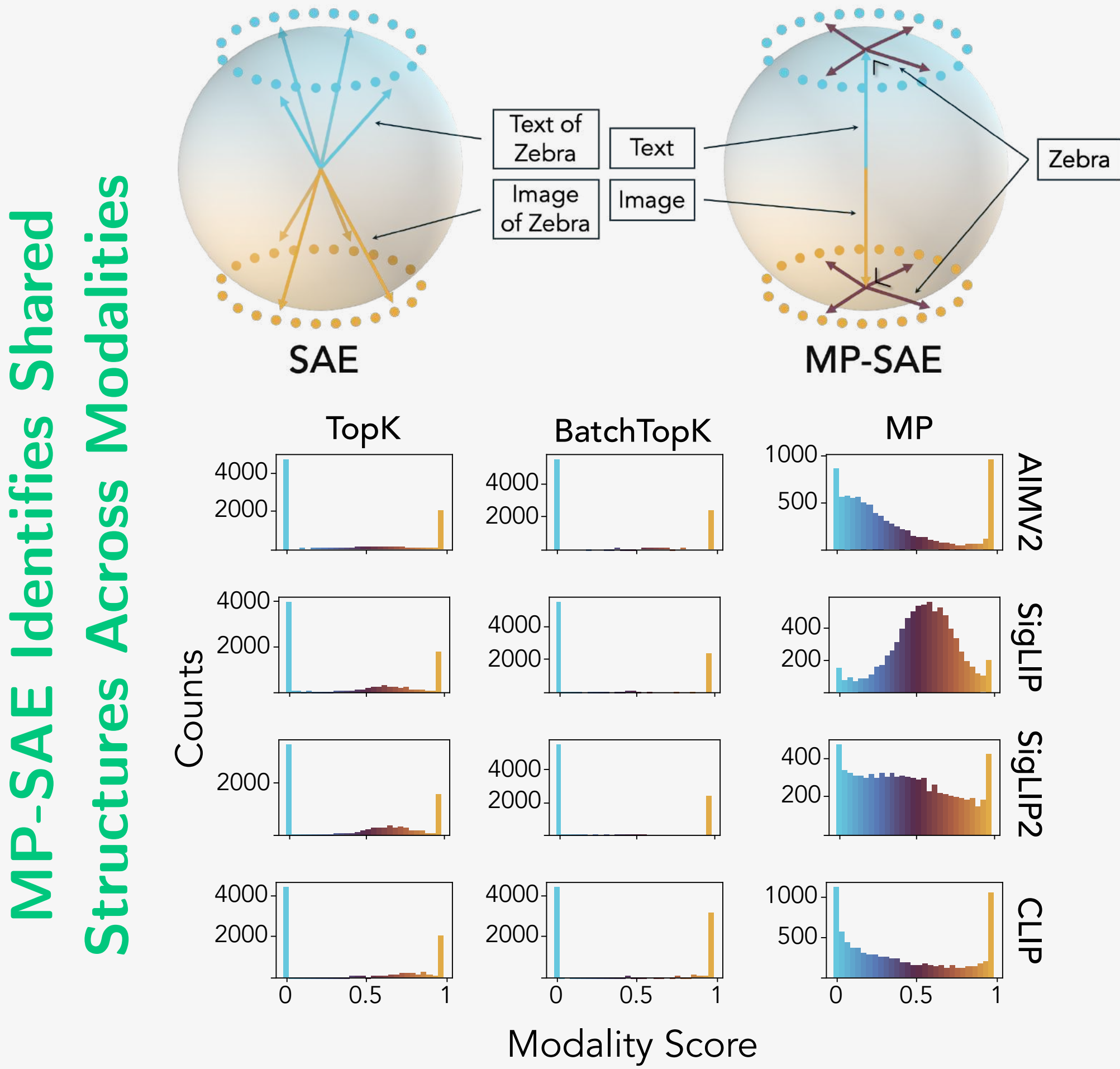
MP-SAE is More Expressive



MP-SAE Learns Correlated Features, Selects Conditionally Orthogonal Ones



Interpreting Multi-Modal Large Models



NeuBahar Lab

NeurIPS 2025

MP-SAE on MNIST

CRISP Group at Harvard