



Shaping Latent Geometry with Noise-Injected Hopfield Dynamics

Wooyul Jung¹, Youngseok Joo¹,

Dohyun Yu², Suhjung Choi¹, Byoung-Tak Zhang¹

²College of Liberal Studies, Seoul National University

¹Interdisciplinary Program in Artificial Intelligence, Seoul National University



Introduction

Motivation.

Representation learning often incidentally acquires geometric structure e.g., class-mean collapse late in training or consistent eigenspectrum decay across biological and artificial systems. Such phenomena suggest that useful latent geometry emerges spontaneously but is not explicitly controlled.

Prior work explicitly shapes latent spaces by (i) controlling local sensitivity and (ii) contrastive objectives, and (iii) explicit priors. Yet these approaches do not exploit attractor dynamics as a direct mechanism for organizing representations.

Approach.

We introduce noise-injected Hopfield retrieval layers inserted between network layers to explicitly sculpt latent spaces. By tuning two parameters # of retrieval iterations, and inverse temperature, we control how strongly representations contract, interpolating from discrete attractors to smooth continuous manifolds.

Effect.

These layers organize intermediate features into stable clusters, naturally improving noise robustness and overall representation quality.

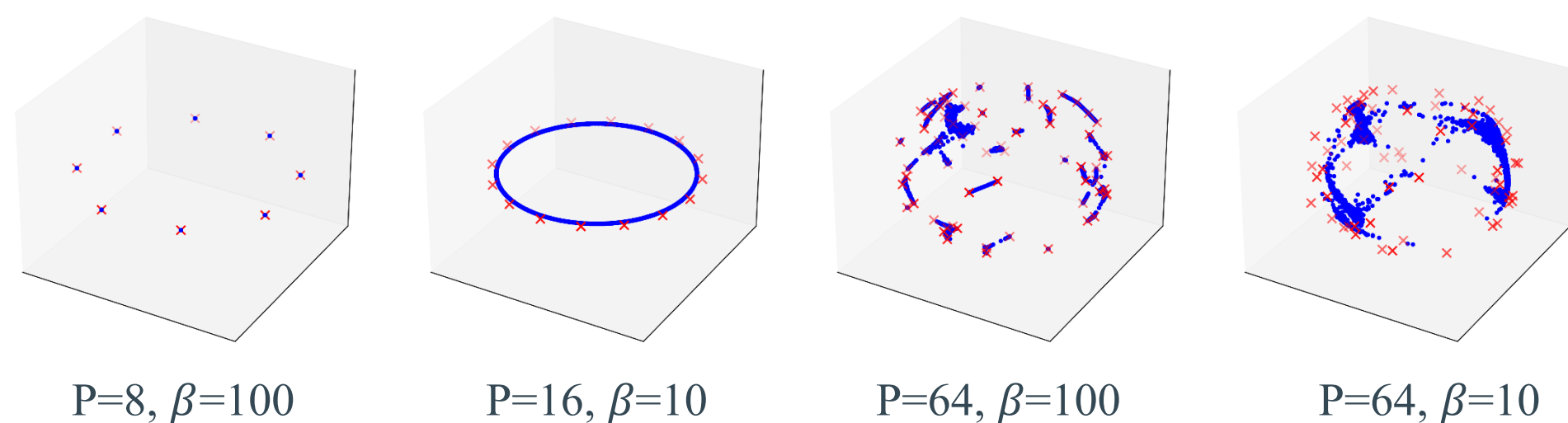
Hopfield Retrieval layer

Modern Hopfield Dynamics.

We formulate a stochastic Hopfield update that introduces noise-driven attractor dynamics into latent representations. Given a set of unit-norm memory vectors $\{e_i\}_{i=1}^P$ with $\|e_i\|_2 = 1$ and memory matrix $C = [e_1, \dots, e_P] \in \mathbb{R}^{D \times P}$, the query vector $z_t \in \mathbb{R}^D$ is updated as follows. At each iteration t , Gaussian noise $\xi_t \sim \mathcal{N}(0, \sigma^2 I_D)$ is injected and the perturbed vector is normalized to unit norm. This process corresponds to gradient descent on the Modern Hopfield energy function.

$$z_{t+1} = C \operatorname{softmax}(\beta C^\top \hat{z}_t), \quad \hat{z}_t = \frac{z_t + \xi_t}{\|z_t + \xi_t\|_2}.$$

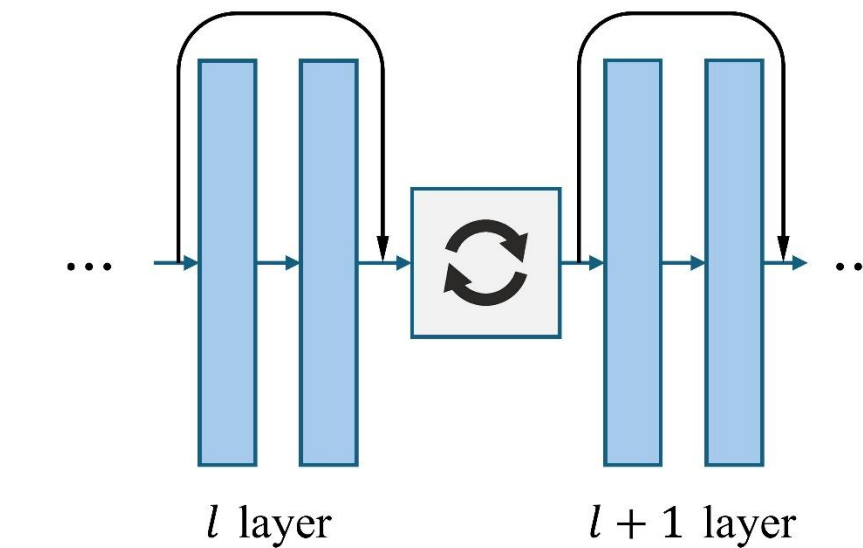
$$E_{\text{MHN}}(z) = -\frac{1}{\beta} \log \sum_{i=1}^P \exp(\beta z^\top e_i).$$



Noise-Injected Hopfield Retrieval (NIHR).

We implement these dynamics as NIHR, a general and fully differentiable module that can be inserted between arbitrary network layers. NIHR performs iterative refinement, pulling representations toward Hopfield energy minima and thereby enforcing local geometric structure in the latent space.

$$z_{t+1} = \mathcal{H}_\sigma(z_t), \quad t = 0, \dots, T-1.$$

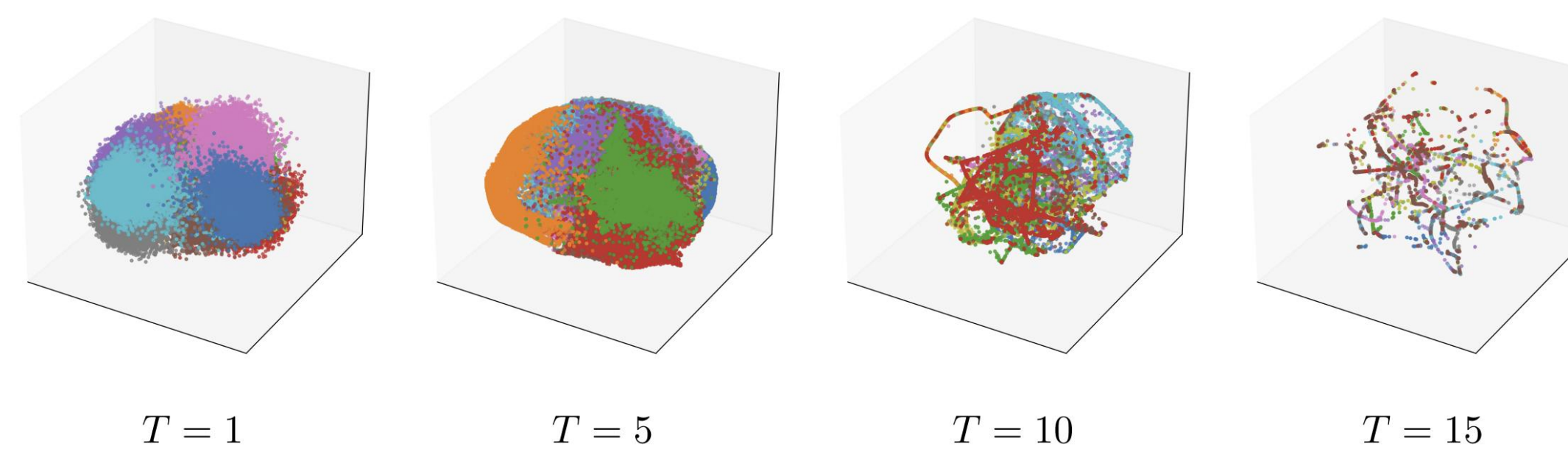


Structural evolution under Autoencoder framework

We insert NIHR between the encoder and decoder to study how Hopfield retrieval reshapes latent organization. This allows us to analyze how structured refinement influences reconstruction robustness and representation quality.

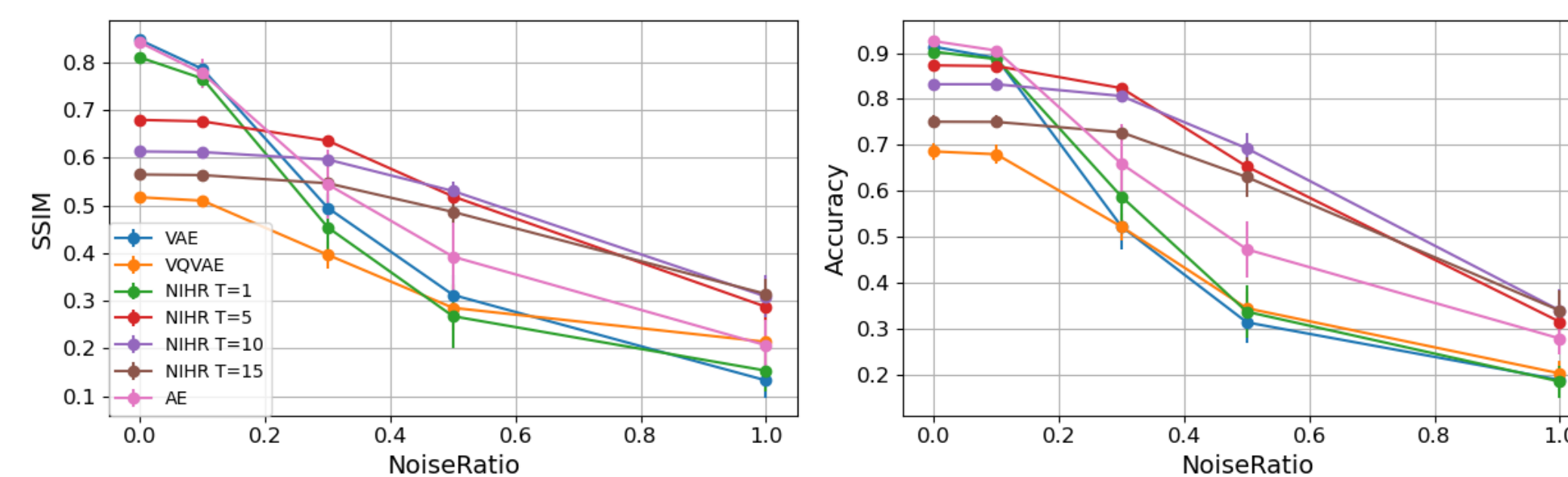
PCA Visualization.

we project latent codes onto their top principal components after $T = 1, 5, 10, 15$ refinement steps. As T increases, diffuse latents progressively contract into clear, stable clusters, demonstrating how NIHR sharpens latent geometry over iterations.



Robustness under Gaussian Noise.

We conduct controlled experiments on MNIST with additive Gaussian noise of varying standard deviation σ to simulate input degradation. Two metrics are used: the structural similarity index (SSIM) for reconstruction quality, and linear probing accuracy for assessing the discriminative power of learned representations.



Evaluation under Gaussian input corruption. Left: SSIM scores for reconstruction fidelity. Right: Linear probing accuracy for representation quality. Models include AE, VAE, VQ-VAE, and NIHR with $T \in \{1, 5, 10, 15\}$.

Robust Classification under Corruptions

We integrate NIHR into ResNet-18 as a latent refinement layer to improve robustness under input corruptions. The insertion point is chosen empirically: after the first residual block for MNIST and after the last residual block for CIFAR-10.

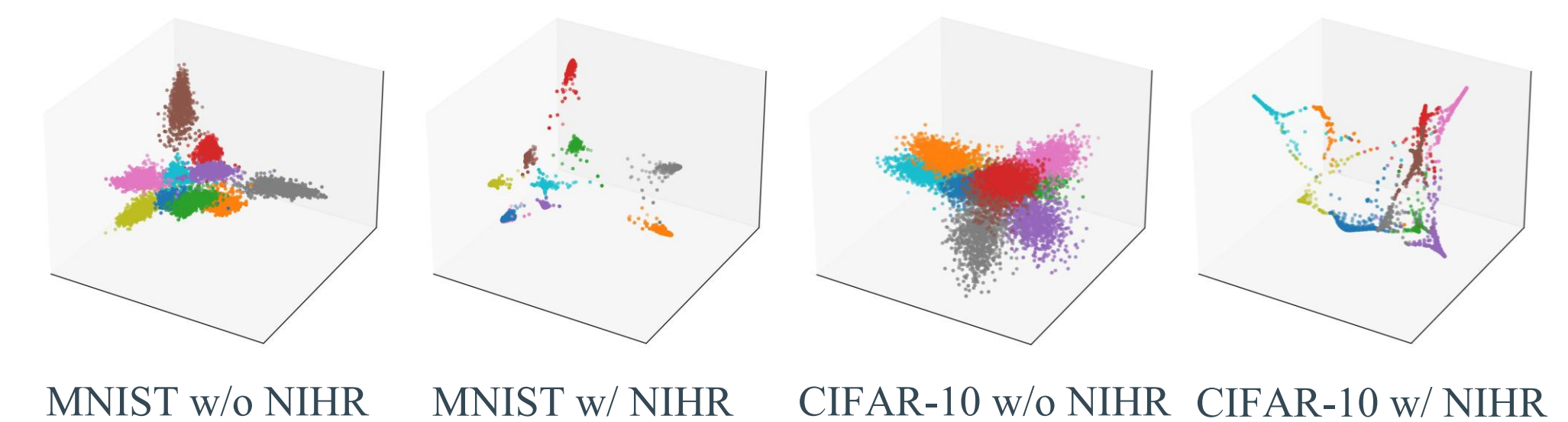
Evaluation under Corruption.

Robustness is measured on MNIST-C and CIFAR-10-C, while all models are trained only on clean data with no augmentation or corruptions. The resulting classification accuracies on corrupted datasets are summarized below.

Dataset	Model	Clean (0)	s1	s2	s3	s4	s5	Overall
MNIST-C	ResNet	99.47 \pm 0.05	–	–	–	–	–	72.13 \pm 2.19
	ResNet+NIHR	99.50 \pm 0.02	–	–	–	–	–	80.56 \pm 2.2
CIFAR-C	ResNet	84.78 \pm 0.18	75.69 \pm 0.68	70.39 \pm 0.92	64.99 \pm 1.03	58.51 \pm 1.21	48.98 \pm 1.04	63.71 \pm 0.96
	ResNet+NIHR	81.81 \pm 0.64	73.93 \pm 1.25	69.28 \pm 1.29	64.58 \pm 1.38	58.73 \pm 1.38	51.00 \pm 1.94	63.50 \pm 1.31
CIFAR-C (Gaussian)	ResNet	84.78 \pm 0.18	70.10 \pm 1.74	55.40 \pm 2.93	41.91 \pm 4.14	36.36 \pm 4.63	31.93 \pm 4.86	47.14 \pm 3.53
	ResNet+NIHR	81.81 \pm 0.64	72.81 \pm 2.30	63.83 \pm 3.41	54.47 \pm 4.21	49.75 \pm 4.59	45.99 \pm 4.75	57.37 \pm 3.63

PCA Visualization.

We visualized how the intermediate representations of ResNet-18 evolve under latent attractor dynamics and observed that they become increasingly coherent as they are pulled toward the attractor.



Contributions

NIHR is a general, differentiable module that integrates into diverse architectures without modification. Its noise-driven Hopfield retrieval forms clear attractor basins that organize intermediate representations, improving robustness to noise and input corruptions. We also plan to extend this work to RNNs and explore the connections among Hopfield networks, Transformers, RNNs, and other ANNs in terms of attractor dynamics.

References

- Arash Vahdat, Karsten Kreis, and Jan Kautz (2021). "Score-based generative modeling in latent space." In: Advances in Neural Information Processing Systems 34, pp. 11287–11302.
- Hubert Ramsauer et al. (2020). "Hopfield networks is all you need." In: arXiv preprint arXiv:2008.02217.

