A BIOLOGICALLY PLAUSIBLE DENSE ASSOCIATIVE MEMORY WITH EXPONENTIAL CAPACITY

Anonymous authorsPaper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

034

037

038

040

041

042

043

044

046

047

048

050 051

052

ABSTRACT

Krotov and Hopfield (2021) proposed a biologically plausible two-layer associative memory network with memory storage capacity exponential in the number of visible neurons. However, the capacity was only linear in the number of hidden neurons. This limitation arose from the choice of nonlinearity between the visible and hidden units, which enforced winner-takes-all dynamics in the hidden layer, thereby restricting each hidden unit to encode only a single memory. We overcome this limitation by introducing a novel associative memory network with a threshold nonlinearity that enables distributed representations. In contrast to winner-takes-all dynamics, where each hidden neuron is tied to an entire memory, our network allows hidden neurons to encode basic components shared across many memories. Consequently, complex patterns are represented through combinations of hidden neurons. These representations reduce redundancy and allow many correlated memories to be stored compositionally. Thus, we achieve much higher capacity: exponential in the number of hidden units, provided the number of visible units is sufficiently larger than the number of hidden neurons. Exponential capacity arises because all binary states of the hidden units can become stable memory patterns with an appropriately chosen threshold. Moreover, the distributed hidden representation, which has much lower dimensionality than the visible layer, preserves class-discriminative structure more effectively than the raw visible patterns, supporting efficient nonlinear decoding. These results establish a new regime for associative memory, enabling high-capacity, robust, and scalable architectures consistent with biological constraints.

1 Introduction

Associative memory networks are a class of attractor models in which the system can recall stored memories from their incomplete or noisy versions via recurrent dynamics (Krotov et al. (2025)). In such models, memories are conceptualized as the stable fixed points of the network dynamics. The number of fixed points determines the storage capacity of the network, and significant efforts have been made to construct networks with sufficiently high storage capacity to explain human memory.

The classical Hopfield network, a leading model for associative memory, has a storage capacity that scales linearly with the number of neurons in the network (Hopfield (1982)). Dense associative memories, sometimes also referred to as modern Hopfield networks (Krotov and Hopfield (2016)), are promising modifications of the classical Hopfield model. By incorporating higher-order interactions (e.g., interactions that are quadratic, rather than linear, in the input to a neuron), they achieve a storage capacity that scales super-linearly with the number of neurons. There are many possible choices for the energy function in this class of models. For instance, the power interaction vertex leads to the power-law scaling of the capacity (Baldi and Venkatesh (1987); Gardner (1987); Abbott and Arian (1987); Horn and Usher (1988); Chen et al. (1986); Krotov and Hopfield (2016)). More sophisticated shapes of the energy function result in exponential storage capacity in the number of neurons, while maintaining large basins of attraction (Demircigil et al. (2017); Lucibello and Mézard (2024)).

The naive implementation of Dense Associative Memory models, however, requires synaptic interactions that are biologically implausible. In particular, these models rely on nonlinear interactions among synapses; e.g., the drive to a neuron could be a quadratic function of its input. While specific biological mechanisms have been proposed to support these interactions, such as astrocytic processes

enabling four-neuron interactions (Kozachkov et al. (2025)), or dendritic computation Kafraj et al. (2025), these remain limited in scope and dictate strong constraints on the possible shape of the energy landscape. To get around the biological implausibility of these models, Krotov and Hopfield (2021) introduced a two-layer model. In this architecture, the visible neurons correspond to features of the patterns, while the hidden neurons serve as auxiliary computational elements that mediate complex interactions. Higher-order interactions among visible neurons emerge by selecting appropriate activation functions for the hidden neurons.

Nevertheless, the two-layer implementation of Krotov and Hopfield has two key limitations. First, the storage capacity is at most linear in the number of hidden neurons (Krotov (2021); Krotov and Hopfield (2021)). This is unsatisfactory from the perspective of information storage – one would like to store as much information as possible while utilizing only a small number of neurons. Second, at inference time, the network demonstrates a winner-takes-all behavior. This means that the asymptotic fixed point that the network converges to corresponds to a single hidden neuron being activated, while the rest of the hidden neurons are inactive. This behavior results in grandmother-like representations for hidden neurons, as opposed to distributed representations, which are more efficient at storing information.

Our work tackles these two limitations. Specifically, we present a novel implementation of Dense Associative Memory that achieves exponential storage capacity in the number of hidden neurons. This is accomplished with a simple yet critical change: we use a threshold activation function that does not enforce winner-takes-all dynamics. The threshold activation enables distributed memory representations—multiple hidden neurons can be active for a memory, and each hidden neuron can participate in multiple memories. As a result, all possible binary patterns of hidden neuron states become stable fixed points, enabling the network to store exponentially many memories, including highly correlated ones. Beyond high capacity, the hidden layer of the network is low-dimensional compared to the visible layer, yet it produces structured representations that preserve class-discriminative information, with memories sharing components represented close together in the hidden activity space. We establish this result through both theoretical analysis and numerical simulations, and show that the resulting fixed points also possess large basins of attraction.

Our model is closely related to the framework recently proposed by Chandra et al. (2025), which combines multiple Dense Associative Memory modules to produce a distributed code for the visible neurons. Each module performs a winner-takes-all operation similar to Krotov and Hopfield (2021), so only a single hidden neuron is active per module. By combining several modules, though, they achieve exponential storage capacity. However, we show that multiple modules are unnecessary: exponential capacity can be achieved with a single module, provided the activation function is chosen appropriately.

Beyond its biological motivation, our work also connects to a growing body of research on Dense Associative Memories in machine learning. Notably, it has been shown that Dense Associative Memory closely corresponds to the attention mechanism in transformer architectures Ramsauer et al. (2021); Hoover et al. (2023a), offering a principled framework for viewing the transformer's attention and feedforward computations as steps in a global energy minimization process. Complementary research has demonstrated that generative diffusion models, widely used in high-quality image generation, also exhibit associative memory behavior Hoover et al. (2023b); Ambrogioni (2024); Pham et al. (2025). Further studies have expanded the model's functionality: for instance, Chaudhry et al. (2023) examined its ability to store and retrieve long sequences; Burns and Fukai (2023) introduced higher-order simplicial interactions; and Dohmatob (2023); Hoover et al. (2025) proposed alternative energy functions that also support exponential storage capacity. Our results contribute to this line of work by showing how exponential storage capacity can be achieved within a biologically plausible two-layer framework, thereby bridging theoretical neuroscience with modern machine learning architectures.

In the following sections, we first formally define the model and its dynamics and derive the optimal threshold analytically for a network with fixed weights. We then present a theoretical analysis of storage capacity and basins of attraction, showing that the network exhibits large basins of attraction, making recall robust to substantial noise in the visible inputs. Next, we introduce a learning rule for storing real, correlated memories, enabling compositional memory storage, and present numerical experiments on MNIST and CIFAR-10 that demonstrate high-capacity recall, structured hidden representations, and robustness to noise. Finally, we conclude by discussing the biological plausibility

of the network and potential directions for extending the model to incorporate additional constraints and more realistic neuronal properties.

2 RESULTS

2.1 Model

In this section, we present our model and demonstrate that its storage capacity scales exponentially with the number of hidden neurons, meaning that all possible binary patterns of hidden neurons are stable fixed points.

We first define the dynamics of the system as follows:

$$\tau_v \frac{dv_i}{dt} = -v_i + \sum_{\mu=1}^{N_h} \tilde{\xi}_{i\mu} \Theta(h_\mu - \theta)$$
 (1a)

$$\tau_h \frac{dh_\mu}{dt} = -h_\mu + \sum_{i=1}^{N_v} \tilde{\xi}_{\mu i} v_i,$$
(1b)

where $\Theta(\cdot)$ is the standard Heaviside step function:

$$\Theta(z) = \begin{cases} 0 & \text{if } z \le 0\\ 1 & \text{if } z > 0 \end{cases}$$
 (2)

The parameter θ will be chosen to ensure the stability of all binary patterns in the hidden layer.

The network consists of N_v visible neurons (the v_i) and N_h hidden neurons (the h_μ), arranged in a bipartite architecture—i.e., without lateral connections within either layer.

For the following analysis of capacity and basins of attraction, we chose the scaling factors of the synaptic connections between a visible neuron i and a hidden neuron μ as:

$$\tilde{\xi}_{i\mu} = \frac{1}{\sqrt{N_h}} \xi_{i\mu} \tag{3a}$$

$$\tilde{\xi}_{\mu i} = \frac{\sqrt{N_h}}{N_n} \xi_{\mu i} \tag{3b}$$

purely for convenience, as they simplify subsequent expressions. These connections are reciprocal and randomly drawn from a standard normal distribution:

$$\xi_{\mu i} = \xi_{i\mu} \sim \mathcal{N}(0, 1). \tag{4}$$

2.2 STORAGE CAPACITY

To determine the storage capacity, we'll first focus on the fixed points of the dynamics given in Eq. (1). Defining

$$s_{\mu} \equiv \Theta(h_{\mu} - \theta) \,, \tag{5}$$

it is straightforward to show that in steady state, s_μ satisfies

$$s_{\mu} = \Theta\left(\sum_{\nu=1}^{N_h} J_{\mu\nu} s_{\nu} - \theta\right) \tag{6}$$

where

$$J_{\mu\nu} \equiv \frac{1}{N_v} \sum_{i=1}^{N_v} \xi_{\mu i} \xi_{i\nu} \,. \tag{7}$$

Equation (6), with the weight matrix given in Eq. (7), is very close to the classical Hopfield model; the only difference is that in the classical model, the $\xi_{\mu i}$ are binary, whereas in our model they're Gaussian. However, the classical Hopfield model works in the regime $N_v < N_h$, with memory storage possible only if $N_v < 0.138N_h$ (Amit et al. (1985)). Here, though, we'll consider a very different regime: $N_v \gg N_h$. In this limit, $J_{\mu\nu}$ approaches the identity matrix (Marchenko and Pastur (1967)), which completely decouples the hidden neurons. Assuming the threshold, θ , is chosen correctly, this leads immediately to exponential storage capacity.

Exponential capacity clearly holds in the limit $N_v \to \infty$. What happens when N_v is finite? We show in Appendix A.1 that

$$J_{\mu\nu} = \delta_{\mu\nu} + \frac{\zeta_{\mu\nu}}{\sqrt{N_r}} \tag{8}$$

where the $\zeta_{\mu\nu}$ are independent, zero-mean, unit-variance Gaussian random variable,

$$\zeta_{\mu\nu} \sim \mathcal{N}(0,1) \,,$$
(9)

and here and in what follows $\delta_{\mu\nu}$ is the Kronecker delta. Thus, Eq. (6) may be written

$$s_{\mu} = \Theta\left(s_{\mu} + \frac{1}{\sqrt{N_v}} \sum_{\nu=1}^{N_h} \zeta_{\mu\nu} s_{\nu} - \theta\right). \tag{10}$$

Because the $\zeta_{\mu\nu}$ are independent, the second term in parentheses, q_{μ} , scales as

$$\left| \frac{1}{\sqrt{N_v}} \sum_{\nu=1}^{N_h} \zeta_{\mu\nu} s_{\nu} \right| \sim \sqrt{\frac{1}{N_v} \sum_{\nu=1}^{N_h} s_{\nu}^2} \le \sqrt{\frac{N_h}{N_v}}$$
 (11)

where the second inequality follows because s_{ν} is either 0 or 1.

If we set $\theta=1/2$, in the limit $N_v\gg N_h$ Eq. (10) typically has two solutions: one at $s_\mu=0$ and one at $s_\mu=1$. In fact, the probability that there is only one solution is the probability that $|q_\mu|>1/2$, which scales at most as $e^{-N_v/2N_h}$. Thus, even when N_v is only about ten times larger than N_h , and the threshold is not exactly 1/2, there are approximately 2^{N_h} fixed points. And if we consider fixed points with, say, at most N of the s_μ nonzero, then we only need N_v on the order of 10N, which can be relatively small.

There are exponentially many fixed points, but are they stable? To answer that, we need to determine the value of h_{μ} at the fixed points. combining Eq. (1) with the definitions of s_{μ} , Eq. (5), and $J_{\mu\nu}$, Eq. (7), we have

$$h_{\mu} = \sum_{\nu} J_{\mu\nu} s_{\nu} . \tag{12}$$

Since $J_{\mu\nu}$ is approximately the identity matrix, we see that at equilibrium h_{μ} is close to either 0 or 1. Thus, because our nonlinearity is a step function, its derivative vanishes at equilibrium, guaranteeing the stability of the fixed points (Appendix A.2). Consequently, when we solve Eq. (1), we expect to see 2^{N_h} stable fixed points in the regime $N_v \gg N_h$. This prediction is consistent with numerical simulations, as can be seen in Figure 1a.

2.3 Basins of Attraction

Although the fixed points are stable, that still leaves the question: how big are the basins of attraction? Since the hidden units have no structure, we'll assume that noisy input enters the network via the visible units, and initially all the h_{μ} are zero. How far from the fixed points can the input be and still be recalled perfectly?

Assuming the noise is additive, the initial values of the visible and hidden neurons are,

$$v_i(0) = \frac{1}{\sqrt{N_h}} \sum_{\nu=1}^{N_h} \xi_{i\mu} \Theta(h_{\mu, \text{target}} - \theta) + \epsilon_i^{\nu}$$
(13a)

$$h_{\mu}(0) = 0 \tag{13b}$$

where $h_{\mu,\text{target}} = 1$ if neuron μ encodes the target memory, and 0 otherwise (motivated by the fact that h_{μ} is close to either 0 or 1 at the fixed points; see Eq. (12)).

To reach the target fixed point in both the hidden and visible unit space, the hidden neurons must evolve to their target values, $h_{\mu, \text{target}}$, before the visible units change much. That requires the visible units to evolve much more slowly than the hidden units, which we can guarantee by setting $\tau_h \ll \tau_v$. With this condition, at a time t satisfying $\tau_h \ll t \ll \tau_v$, $h_\mu(t)$ reaches equilibrium while $v_i(t)$ is still approximately equal to $v_i(0)$. Using Eq. (1a) with $dv_i/dt=0$ along with Eq. (7), that equilibrium is given by

$$h_{\mu}(t) = \sum_{\nu} J_{\mu\nu} \Theta(h_{\nu,\text{target}} - \theta) + \frac{\sqrt{N_h}}{N_v} \sum_{i=1}^{N_v} \xi_{\mu i} \epsilon_i^v + \mathcal{O}(t/\tau_v).$$
 (14)

Using Eq. (8), we see that the first term is $\Theta(h_{\mu,\text{target}} - \theta) + \mathcal{O}(\sqrt{N_h/N_v})$. And the second term scales as $\sigma_v \sqrt{N_h/N_v}$ where σ_v^2 is the variance of the noise. Thus, so long as

$$\operatorname{Var}[\epsilon] \ll \frac{N_v}{N_h},\tag{15}$$

 $h_{\mu}(t)$ will be close to its target value when $t \ll \tau_v$. Since $v_i(t)$ is close to its target value at that time, it will stay close, and asymptotically the target pattern will be recovered. Given that $N_v \gg N_h$, ϵ_i^v can be very large without affecting recall. Thus, the basin of attraction is very large (see Figure 1b).

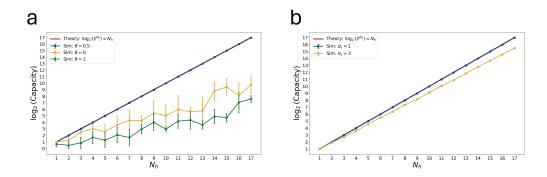


Figure 1: Capacity versus the number of hidden units, N_h , with $N_v = 100N_h$ and $\tau_v = 20\tau_h$. (a) Capacity for different thresholds, θ . The highest storage capacity is achieved when the threshold is set to its optimal theoretical value, $\theta = 0.5$. (b) The effect of noise in the visible layer (ϵ_i^v in Eq. (13a)), shown for different noise variances, demonstrates the large basin of attraction of the fixed points.

2.4 Learning Rule

So far we have focused on storage capacity with fixed synaptic weights. A natural next step is to understand how these weights can be learned. In this section, we introduce a learning rule that reflects *compositional learning*: a small number of simple, reusable components can be combined to form complex patterns, and conversely, complex patterns can be decomposed into simpler components.

Assuming symmetric weights, $\tilde{\xi}_{i\mu}=\tilde{\xi}_{\mu i}$, the steady-state visible activity in Eq. (1) can be expressed as

$$\mathbf{v} = \sum_{\mu=1}^{N_h} \tilde{\boldsymbol{\xi}}_{\mu} s_{\mu},\tag{16}$$

where $\tilde{\xi}_{\mu}$ is the μ -th column of $\tilde{\xi} \in \mathbb{R}^{N_v \times N_h}$, i.e., $(\tilde{\xi}_{\mu})_i = \tilde{\xi}_{i\mu}$. If only hidden neuron μ is active, the visible state equals $\tilde{\xi}_{\mu}$. A visible memory is thus called *basic* if it corresponds to a single active

hidden neuron, and *complex* if it is formed by the activation of multiple hidden neurons, i.e. a composition of several basic memories.

The goal of learning is to find a synaptic weight matrix $\tilde{\xi}$ and a threshold θ such that a set of target memories $\{\mathbf{v}_m \in \mathbb{R}^{N_v}\}_{m=1}^M$ approximately correspond to stable fixed points of the network dynamics, with $M \gg N_h$ (e.g., MNIST or CIFAR-10). This is achieved using the following optimization procedure,

$$(\tilde{\xi}, \theta) = \arg\min_{\tilde{\xi}, \theta} \sum_{m=1}^{M} \left\| \mathbf{v}_m - \sum_{\mu=1}^{N_h} \tilde{\xi}_{\mu} \Theta(\tilde{\xi}_{\mu}^{\top} \mathbf{v}_m - \theta) \right\|^2,$$
(17)

where s_{μ} is replaced by its target steady-state value. This learning rule is identical to the one proposed in Radhakrishnan et al. (2020). We used Xavier initialization for the weights and approximated the threshold function Θ with a sharp sigmoid to allow gradient-based training.

Figure 2 shows recall results after training on 60,000 MNIST digits with $N_v = 784$ and $N_h = 50$. Despite the high correlation among patterns, the network learns 55,376 unique minima. Variants of the same digit produce hidden representations that are distinct yet partially overlapping, and the recalled visible states remain recognizable.

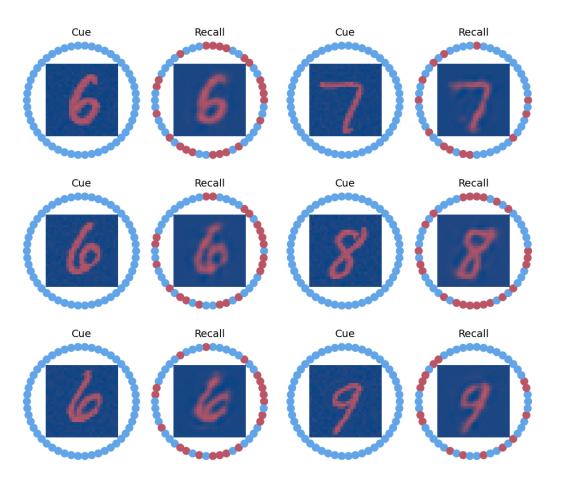


Figure 2: Examples of recall in a network with 50 hidden neurons that memorized 60,000 MNIST images. Hidden neurons are shown on the ring, and visible neurons are visualized as two-dimensional images. Red indicates high activity; blue indicates low activity. Highly correlated images of every digit, for instance, the digit 6 shown here, converge to unique but overlapping hidden representations.

Figure 3a shows the learned basic memories for the MNIST dataset, which correspond to the columns of $\tilde{\xi}$. As shown in Figure 3b, these basic memories are nearly orthogonal, consistent with Eq. (8).

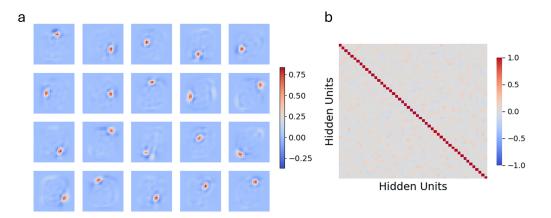


Figure 3: a) 20 (out of 50) columns of the learned weight matrix, which serve as basic memories, are shown as two-dimensional images. b) Correlation matrix of the basic patterns, which correspond to the hidden units.

To evaluate the generality of the proposed learning rule beyond MNIST, we applied the same procedure to grayscale-normalized ([0, 1]) images from the CIFAR-10 dataset. In this case, $N_v=1024$, and to compensate for the increased complexity of this dataset, we used a network with 500 hidden neurons, compared to 50 hidden neurons for MNIST. Figure 4 presents examples of the cues alongside their recalls. These results show that the network is able to reconstruct interpretable outputs from the learned representations, despite storing a large number of complex memories (50,000) and significantly violating the condition $N_v\gg N_h$. Importantly, these images are highly correlated, yet the network produces 49,988 unique stable minima, with each memory representation being both stable and interpretable. The learned basic memories for CIFAR-10 images are shown in Figure 5a. They form a more heterogeneous set, yet remain nearly orthogonal, as shown in Figure 5b.

Finally, to evaluate how well the network preserves class-discriminative information, classifiers were trained on the recalled hidden and visible representations as well as on the original images, using both MNIST and CIFAR-10. Both a linear classifier (logistic regression) and a nonlinear multilayer perceptron (MLP) with two hidden layers were used. The hidden and visible representations preserve class information very effectively compared to the original images. Importantly, the lower-dimensional hidden layer retains this information almost perfectly, demonstrating that its encoding is structured and meaningful: correlated memories are represented close together and remain classifiable.

| Representation | MNIST Accuracy | | CIFAR-10 Accuracy | |
|---------------------------|----------------|-----------|-------------------|-----------|
| | Linear | Nonlinear | Linear | Nonlinear |
| Recalled Hidden Patterns | 86% | 100% | 35% | 99% |
| Recalled Visible Patterns | 91% | 100% | 34% | 70% |
| Original Images | 94% | 100% | 34% | 66% |

Table 1: Classification accuracy of linear (logistic regression) and nonlinear (MLP) classifiers on recalled hidden and visible representations and original images for MNIST and CIFAR-10. Both representations preserve class-discriminative information very well compared to the original images.

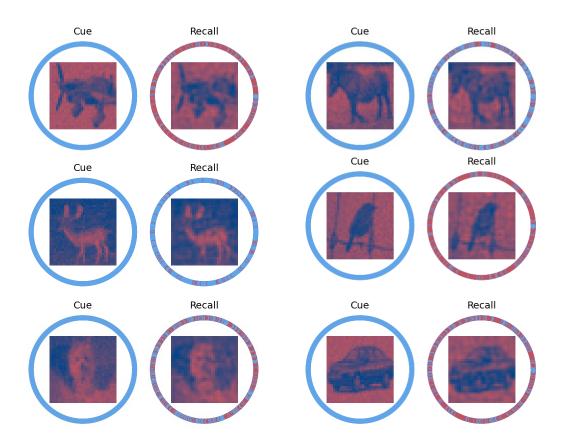


Figure 4: Examples of recall in a network with 500 hidden neurons that memorized 50,000 CIFAR-10 images. Hidden neurons are arranged on a ring, while visible neurons are shown as two-dimensional images. Red indicates high activity, and blue indicates low activity.

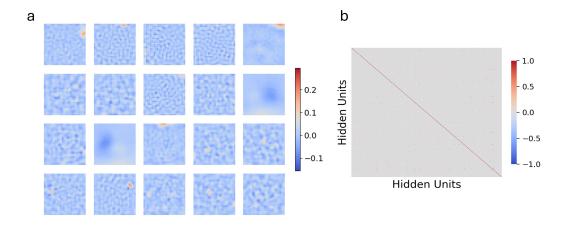


Figure 5: 20 (out of 500) learned basic patterns from grayscale CIFAR-10. b) Orthogonality of these basic patterns, where 0 indicates complete orthogonality. c) Correlation matrix of the basic patterns, which correspond to the hidden units.

3 CONCLUSION

This work introduces a novel Dense Associative Memory Krotov and Hopfield (2021) that achieves exponential storage capacity in the number of hidden neurons, overcoming the limitations of previous two-layer models. By using a threshold activation function, with a theoretically derived threshold, the network supports distributed hidden representations, allowing each hidden neuron to participate in multiple memories. This enables compositional storage of complex and correlated patterns, reducing redundancy while maintaining robust retrieval.

Specifically, the network achieves exponential capacity, 2^{N_h} , using only $N_h N_v$ parameters. In contrast, previous two-layer implementations were limited to a maximum capacity of N_h (Krotov and Hopfield (2021)). As a result, the number of memories per weight grows as $\frac{2^{N_h}}{N_h N_v} \approx 2^{N_h}$, while in previous implementations it is at best $\frac{1}{N_v}$. Even for complex datasets such as MNIST and CIFAR-10, networks with only 50 and 500 hidden units, respectively, were able to store tens of thousands of highly correlated memories and associate the vast majority of them with unique minima, whereas previous models could not store more memories than the number of hidden units.

The model is biologically grounded, relying solely on standard pairwise synapses, and its fixed points have large basins of attraction, ensuring robust recall from noisy inputs. Moreover, the hidden layer produces low-dimensional representations that preserve class-discriminative information, organizing memories with shared components close together in the activity space. This structured representation supports efficient nonlinear decoding that outperforms the raw visible patterns.

Overall, this work establishes a new regime for associative memory, combining high capacity, robust recall, compositional and interpretable representations, and biological plausibility. It provides a theoretical foundation for scalable memory systems that bridge neuroscience models and modern machine learning architectures. Future work will explore the model's capacity under additional biological constraints, including sparse connectivity, compliance with Dale's law, and realistic, unsaturated neuronal firing rates.

REFERENCES

- Laurence F Abbott and Yair Arian. Storage capacity of generalized networks. *Physical review A*, 36 (10):5091, 1987.
- Luca Ambrogioni. In search of dispersed memories: Generative diffusion models are associative memory networks. *Entropy*, 26(5):381, 2024.
 - Daniel J Amit, Hanoch Gutfreund, and Haim Sompolinsky. Storing infinite numbers of patterns in a spin-glass model of neural networks. *Physical Review Letters*, 55(14):1530, 1985.
 - Pierre Baldi and Santosh S Venkatesh. Number of stable points for spin-glasses and neural networks of higher orders. *Physical Review Letters*, 58(9):913, 1987.
 - Thomas F Burns and Tomoki Fukai. Simplicial hopfield networks. In *The Eleventh International Conference on Learning Representations*, 2023.
 - Sarthak Chandra, Sugandha Sharma, Rishidev Chaudhuri, and Ila Fiete. Episodic and associative memory from spatial scaffolds in the hippocampus. *Nature*, 2025.
 - Hamza Chaudhry, Jacob Zavatone-Veth, Dmitry Krotov, and Cengiz Pehlevan. Long sequence hopfield memory. *Advances in Neural Information Processing Systems*, 36:54300–54340, 2023.
 - HH Chen, YC Lee, GZ Sun, HY Lee, Tom Maxwell, and C Lee Giles. High order correlation model for associative memory. In *AIP Conference Proceedings*, volume 151, pages 86–99, 1986.
 - Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Upgang, and Franck Vermet. On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168, 2017.
 - Elvis Dohmatob. A different route to exponential storage capacity. In *Associative Memory* {\&} *Hopfield Networks in 2023*, 2023.
 - Elizabeth Gardner. Multiconnected neural network models. *Journal of Physics A: Mathematical and General*, 20(11):3453, 1987.
 - Benjamin Hoover, Yuchen Liang, Bao Pham, Rameswar Panda, Hendrik Strobelt, Duen Horng Chau, Mohammed Zaki, and Dmitry Krotov. Energy transformer. *Advances in neural information processing systems*, 36:27532–27559, 2023a.
 - Benjamin Hoover, Hendrik Strobelt, Dmitry Krotov, Judy Hoffman, Zsolt Kira, and Duen Horng Chau. Memory in plain sight: A survey of the uncanny resemblances between diffusion models and associative memories. In *Associative Memory* {\&} *Hopfield Networks in 2023*, 2023b.
 - Benjamin Hoover, Zhaoyang Shi, Krishnakumar Balasubramanian, Dmitry Krotov, and Parikshit Ram. Dense associative memory with epanechnikov energy. *arXiv preprint arXiv:2506.10801*, 2025.
 - John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79:2554–2558, 1982.
 - D Horn and M Usher. Capacities of multiconnected memory models. *Journal de Physique*, 49(3): 389–395, 1988.
 - Mohadeseh Shafiei Kafraj, Dmitry Krotov, Brendan A Bicknell, and Peter E Latham. A biologically plausible associative memory network. In *New Frontiers in Associative Memories*, 2025.
- Leo Kozachkov, Jean-Jacques Slotine, and Dmitry Krotov. Neuron-astrocyte associative memory.

 Proceedings of the National Academy of Sciences of the United States of America, 122(21): e2417788122, 2025.
 - Dmitry Krotov. Hierarchical associative memory. arXiv preprint arXiv:2107.06446, 2021.
 - Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29, 2016.

Dmitry Krotov and John J Hopfield. Large associative memory problem in neurobiology and machine learning. In International Conference on Learning Representations, 2021. Dmitry Krotov, Benjamin Hoover, Parikshit Ram, and Bao Pham. Modern methods in associative memory. arXiv preprint arXiv:2507.06211, 2025. Carlo Lucibello and Marc Mézard. Exponential capacity of dense associative memories. Physical Review Letters, 132(7):077301, 2024. VA Marchenko and Leonid A Pastur. Distribution of eigenvalues for some sets of random matrices. *Mat. Sb.(NS)*, 72(114):4, 1967. Bao Pham, Gabriel Raya, Matteo Negri, Mohammed J Zaki, Luca Ambrogioni, and Dmitry Krotov. Memorization to generalization: Emergence of diffusion models from associative memory. arXiv preprint arXiv:2505.21777, 2025. Adityanarayanan Radhakrishnan, Mikhail Belkin, and Caroline Uhler. Overparameterized neural networks implement associative memory. Proceedings of the National Academy of Sciences, 117, 2020. Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas Gruber, Markus Holzleitner, Thomas Adler, David Kreil, Michael K Kopp, et al. Hopfield networks is all you need. In International Conference on Learning Representations, 2021.

A TECHNICAL APPENDICES AND SUPPLEMENTARY MATERIAL

A.1 Distributional Properties of $\zeta_{\mu\nu}$

We define the matrix elements

$$\zeta_{\mu\nu} = \frac{1}{\sqrt{N_v}} \sum_{i=1}^{N_v} \xi_{\mu i} \xi_{i\nu}, \quad \mu \neq \nu,$$
(18)

where $\xi_{\mu i}$ s are randomly drawn from a standard normal distribution, Eq. (4).

Each product $\xi_{\mu i}\xi_{i\nu}$ is a zero-mean random variable, since $\xi_{\mu i}$ and $\xi_{i\nu}$ are independent with zero mean.

By the central limit theorem, the sum of these N_v independent terms converges in distribution to a Gaussian. Specifically,

$$\sqrt{N_v} \left(\frac{1}{N_v} \sum_{i=1}^{N_v} \xi_{\mu i} \xi_{i\nu} - \mathbb{E}[\xi_{\mu i} \xi_{i\nu}] \right) \stackrel{d}{\longrightarrow} \mathcal{N}(0, 1), \tag{19}$$

given that $\mathbb{E}[\xi_{\mu i}\xi_{i\nu}]=0$, and $\text{Var}[\xi_{\mu i}\xi_{i\nu}]=1$,

$$\zeta_{\mu\nu} \xrightarrow{d} \mathcal{N}(0,1).$$
 (20)

Now consider the random variable x_{μ} defined as:

$$x_{\mu} = \frac{1}{\sqrt{N_v}} \sum_{\nu=1}^{N_h} \zeta_{\mu\nu} s_{\nu},\tag{21}$$

This is a random variable with respect to the index μ with s_{ν} fixed. Its variance is given by

$$\operatorname{Var}\left[\frac{1}{\sqrt{N_v}}\sum_{\nu=1}^{N_h}\zeta_{\mu\nu}s_{\nu}\right] = \frac{1}{N_v}\sum_{\nu=1}^{N_h}s_{\nu}^{2}\operatorname{Var}[\zeta_{\mu\nu}] = \frac{1}{N_v}\sum_{\nu=1}^{N_h}s_{\nu}^{2}$$
(22)

where we used the fact that the $\zeta_{\mu\nu}$ are independent random variables with mean 0 and variance 1.

A.2 STABILITY OF THE FIXED POINTS

The stability of fixed points is determined by the Jacobian of the system. Grouping the variables into (\mathbf{v}, \mathbf{h}) , corresponding to the visible and hidden units respectively, the Jacobian has the block structure

$$\mathbf{A} = egin{bmatrix} \mathbf{A}_{vv} & \mathbf{A}_{vh} \ \mathbf{A}_{hv} & \mathbf{A}_{hh} \end{bmatrix}.$$

For the diagonal blocks, consider first the visible units. We have

$$\frac{\partial \dot{v}_i}{\partial v_j} = \begin{cases} -1, & j=i, \\ 0, & j\neq i, \end{cases} \quad \Rightarrow \quad \mathbf{A}_{vv} = -\mathbf{I}_{N_v},$$

where \mathbf{I}_{N_v} is the $N_v \times N_v$ identity matrix. Similarly, for the hidden units,

$$\frac{\partial \dot{h}_{\mu}}{\partial h_{\nu}} = \begin{cases} -1, & \nu = \mu, \\ 0, & \nu \neq \mu, \end{cases} \Rightarrow \mathbf{A}_{hh} = -\mathbf{I}_{N_h},$$

where \mathbf{I}_{N_h} is the $N_h \times N_h$ identity matrix.

For the off-diagonal blocks, the derivative of the Heaviside step function in Eq. (2) is zero almost everywhere,

$$\Theta'(z) = 0, \qquad z \neq 0.$$

 Therefore, away from threshold crossings $(h_{\mu} \neq \theta)$ in Eq. (1a),

$$\frac{\partial \dot{v}_i}{\partial h_{\mu}} = 0 \quad \Rightarrow \quad \mathbf{A}_{vh} = \mathbf{0}.$$

The hidden dynamics depend explicitly on the visible variables:

$$\frac{\partial \dot{h}_{\mu}}{\partial v_{i}} = \tilde{\xi}_{\mu i}, \quad \Rightarrow \quad \mathbf{A}_{hv} = (\tilde{\xi}_{\mu i}).$$

Putting everything together, the Jacobian is lower-triangular,

$$\mathbf{A} = egin{bmatrix} -\mathbf{I}_{N_v} & \mathbf{0} \\ \mathbf{A}_{hv} & -\mathbf{I}_{N_h} \end{bmatrix}.$$

The eigenvalues of a triangular matrix are its diagonal entries, which in this case are all equal to -1. Hence, all fixed points of the dynamics are stable.