ADAPTIVE HOPFIELD NETWORK: RETHINKING SIMI-LARITIES IN ASSOCIATIVE MEMORY

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

033

037

038

040

041

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Associative memory models are content-addressable memory systems fundamental to biological intelligence and are notable for their high interpretability. However, existing models evaluate the quality of retrieval based on proximity, which cannot guarantee that the retrieved pattern has the strongest association with the query, failing correctness. We reframe this problem by proposing that a query is a generative variant of a stored memory pattern, and define a variant distribution to model this subtle context-dependent generative process. Consequently, correct retrieval should return the memory pattern with the maximum a posteriori probability of being the query's origin. This perspective reveals that an ideal similarity measure should approximate the likelihood of each stored pattern generating the query in accordance with variant distribution, which is impossible for fixed and pre-defined similarities used by existing associative memories. To this end, we develop adaptive similarity, a novel mechanism that learns to approximate this insightful but unknown likelihood from samples drawn from context, aiming for correct retrieval. We theoretically prove that our proposed adaptive similarity achieves optimal correct retrieval under three canonical and widely applicable types of variants: noisy, masked, and biased. We integrate this mechanism into a novel adaptive Hopfield network (A-Hop), and empirical results show that it achieves state-of-the-art performance across diverse tasks, including memory retrieval, tabular classification, image classification, and multiple instance learning. Our code is publicly available here.

1 Introduction

Associative memory represents a fundamental paradigm in information storage and retrieval, functioning as a content-addressable memory system that serves as a cornerstone of biological intelligence (Miyashita, 1988; Pearce & Bouton, 2001), particularly in the hippocampus and neocortex (Wang et al., 2014). Unlike conventional computer memory, which retrieves data based on a specific address, associative memory retrieves stored patterns by using a partial or noisy variant of the pattern itself as a cue. This memory paradigm enables robust pattern completion, error correction, and fault-tolerant information processing, making it a compelling model for both understanding biological cognition and developing artificial intelligence systems.

The computational modeling of associative memory has evolved dramatically since its inception. Hopfield (1982) pioneered this field by introducing a recurrent neural network, dubbed Hopfield network, capable of storing and retrieving patterns through energy minimization. Subsequent work (Krotov & Hopfield, 2016; Demircigil et al., 2017) extended memory capacity using a steeper energy function. A pivotal breakthrough came with the establishment of a profound connection between the modern Hopfield network and the attention mechanism (Vaswani et al., 2017), achieved by using the softmax(\cdot) function to further separate memories (Ramsauer et al., 2021). This insight not only unified two previously disparate fields but also inspired further refinements that strengthened associative memory's performance from different perspectives (Millidge et al., 2022; Hu et al., 2023; Wu et al., 2024a), and broadened applications to tasks like clustering (Saha et al., 2023), time series prediction (Wu et al., 2024b), and more (Krotov et al., 2025).

Despite these significant advances, a critical and unaddressed limitation pervades the literature: the absence of a rigorous framework for assessing retrieval accuracy. Current evaluations typically rely

on proximity-based criteria, such as ϵ -retrieval (Ramsauer et al., 2021; Hu et al., 2023; Wu et al., 2024a; Hu et al., 2024; 2025), which deem retrieval successful if the retrieved pattern is sufficiently close to a certain stored pattern. However, proximity does not establish correctness; ensuring the retrieval is a valid memory provides no guarantee that it is the *correct* one, that is, the one that has the strongest association with the query. This oversight leads to a universal reliance on fixed, pre-defined similarity measures (e.g., inner product or Euclidean distance between two memory patterns). Such one-size-fits-all metrics fail to capture the nuanced, context-dependent *association*, or *similarity*, between the query and the stored memory patterns. For instance, the word *click* is semantically similar to *tap*, phonetically similar to *clique*, and orthographically to *clock* — illustrating that an appropriate notion of similarity is context and task dependent while fixed metrics cannot adapt to such context, nor can they certify correctness.

Our central premise is that correctness is inherently generative: a query \mathbf{x} emerges as a *variant* of an unknown stored pattern $\boldsymbol{\xi}_k$. So, to properly define and achieve correct retrieval, we should model the generative process that transforms a stored pattern $\boldsymbol{\xi}_k$ into a query \mathbf{x} . To this end, we encapsulate the context-dependent and application-related subtleness into a probabilistic framework centered on the concept of *variant distribution* $\mathcal{V}(\boldsymbol{\xi}_{1\cdots N})$, a joint distribution over stored patterns $\boldsymbol{\xi}_{1\cdots N}$ and memory variants \mathbf{x} , where the likelihood $p_{\mathcal{V}}(\boldsymbol{\xi}_k, \mathbf{x})$ captures how probable that we observe $\boldsymbol{\xi}_k$ and it coincidentally generates \mathbf{x} for $(\boldsymbol{\xi}_k, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\xi}_{1\cdots N})$. Under this view, a *correct retrieval* returns the memory pattern $\boldsymbol{\xi}_k$ maximizing the posterior $p_{\mathcal{V}}(\boldsymbol{\xi}_k|\mathbf{x})$, that is, the likelihood of \mathbf{x} originates from $\boldsymbol{\xi}_k$ when observed \mathbf{x} as query. With further decomposition, maximizing $p_{\mathcal{V}}(\boldsymbol{\xi}_k|\mathbf{x})$ is equivalent to maximizing $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_k)$, i.e., given $\boldsymbol{\xi}_k$, how probable would it generates \mathbf{x} as its variant. The correct retrieval is therefore finding the pattern $\boldsymbol{\xi}_k$ that is most likely to produce \mathbf{x} by varianting itself. This perspective yields an insight that optimal correct retrieval can be achieved by forcing the similarity measure to mimic the behavior of $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_k)$.

However, it is not possible to derive the variant distribution $\mathcal{V}(\mathbf{x}_1...N)$ and the likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_k)$ on most occasions. Thus, we need to reconstruct the unknown by mining deeply from what is observable: the query \mathbf{x} , stored patterns $\boldsymbol{\xi}_{1...N}$, and samples matching the context that vaguely describe \mathcal{V} . Building on these motivations, we introduce an *adaptive similarity* framework that learns to approximate $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_k)$ from samples observed from the variant distribution, without assuming the variant type is known a priori. Integrating this novel similarity measure into the Hopfield energy yields an adaptive Hopfield network that strives for correct retrieval by capturing the underlying variant distribution. Our key contributions are as follows:

- We introduce the **variant distribution** to model how queries emerge from stored patterns, and formalize **correct retrieval** as a robust and meaningful criterion for evaluating the theoretical accuracy of associative memories.
- We propose **adaptive similarity** derived from this framework and prove its optimality for three canonical and widely applicable types of memory variants: noisy, masked, and biased.
- We build a novel **adaptive Hopfield network** (A-Hop) that incorporates learnable adaptive similarity, achieving state-of-the-art performance among computational associative memories on tasks including memory retrieval, tabular classification, image classification, and multiple instance learning.

2 BACKGROUND

We consider an associative memory that stores N memory patterns denoted by the memory matrix $\mathbf{\Xi} = [\boldsymbol{\xi}_1; \boldsymbol{\xi}_2; \cdots; \boldsymbol{\xi}_N] \in \mathbb{R}^{d \times N}$, where each column vector $\boldsymbol{\xi}_i \in \mathbb{R}^d$ represents a memory pattern. Given a memory variant (query) $\mathbf{x} \in \mathbb{R}^d$, the goal is to retrieve the stored memory that is most associated with it. For simplicity, we denote $[n] \triangleq \{k \in \mathbb{Z} \mid 1 \leq k \leq n\}$, and $\boldsymbol{\xi} \in \boldsymbol{\Xi}$ means $\boldsymbol{\xi}$ is one of the column vectors of the memory matrix $\boldsymbol{\Xi}$. Appendix $\boldsymbol{A}.1$ contains a collection of notations.

2.1 HOPFIELD NETWORKS

Hopfield network is a line of associative memory that retrieves the most relevant stored memory through a similarity-based matching process. The original Hopfield network (Hopfield, 1982) uses d binary neurons $\sigma \in \{-1, +1\}^d$ to represent the states of the memory system that is limited to

Table 1: Summary of all Hopfield network by components; Hop (Hopfield, 1982), D-Hop (Krotov & Hopfield, 2016), E-Hop (Demircigil et al., 2017), M-Hop (Ramsauer et al., 2021), U-Hop (Millidge et al., 2022), S-Hop (Hu et al., 2023), K-Hop (Wu et al., 2024a), and A-Hop (Ours).

Model	$sim(\boldsymbol{\xi}, \mathbf{x})$	sep(s)	$\operatorname{mod}(\boldsymbol{\xi})$	$ E(\mathbf{x}) $
Hop (Original)	$oldsymbol{\xi}^{ op}\mathbf{x}$	\mathbf{s}	ξ	$-rac{1}{2}\mathbf{x}^{T}\mathbf{\Xi}\mathbf{\Xi}^{T}\mathbf{x}$
D-Hop (Dense)	$oldsymbol{\xi}^{ op}\mathbf{x}$	\mathbf{s}^k	ξ	$-(ilde{1}^{ op}\mathbf{\Xi}^{ op}\mathbf{x})^{k+1}$
E-Hop (Exponential)	${oldsymbol{\xi}}^{ op}\mathbf{x}$	$\exp(\mathbf{s})$	ξ	$-\exp(1^{\! op}\mathbf{\Xi}^{\! op}\mathbf{x})$
M-Hop (Modern)	${oldsymbol{\xi}}^{ op} \mathbf{x}$	softmax(s)	ξ	$\mathbf{x}^{T}\mathbf{x}/2 - \mathrm{lse}(\mathbf{\Xi}^{T}\mathbf{x})$
U-Hop (Universal)	$-\ {m \xi} - {f x}\ _1$	$arg \max(\mathbf{s})$	ξ	/
S-Hop (Sparse)	$\boldsymbol{\xi}^{ op}\mathbf{x}$	sparsemax(s)	ξ	$\mathbf{x}^{T}\mathbf{x}/2 - \Psi_2^{\star}(\beta \mathbf{\Xi}^{T}\mathbf{x})$
K-Hop (Kernelized)	$\boldsymbol{\xi}^{ op}\mathbf{x}$	α -entmax(\mathbf{s})	$\boldsymbol{\Phi}^{\!\top}\!\boldsymbol{\Phi}\boldsymbol{\xi}$	$\mathbf{x}^{T} \mathbf{\Phi}^{T} \mathbf{\Phi} \mathbf{x} / 2 - \Psi_{\alpha}^{\star} (\beta \mathbf{\Xi}^{T} \mathbf{\Phi}^{T} \mathbf{\Phi} \mathbf{x})$
A-Hop (Adaptive)	$\mathbf{w}^{\top}\mathbf{U}\tilde{\mathbf{q}}$	multiple	$\boldsymbol{\xi}$ or $\boldsymbol{\Phi}^{T} \boldsymbol{\Phi} \boldsymbol{\xi}$	$-\mathrm{lse}(\mathbf{s}(\mathbf{\Xi},\mathbf{x}))$

storage of binary values. For retrieving, the model sets the query as the initial state (i.e., $\sigma^{(0)} = x$), and updates one or more neuron(s) iteratively through the following dynamics until convergence:

$$oldsymbol{\sigma}_i^{(t+1)} = ext{sgn}\left(\sum_{j=1}^d \mathbf{T}_{i,j} oldsymbol{\sigma}_j^{(t)}
ight), \qquad ext{where } \mathbf{T}_{i,j} = \sum_{k=1}^N oldsymbol{\xi}_{k,i} \cdot oldsymbol{\xi}_{k,j}.$$

A vectorized retrieval dynamics exists when updating all neurons simultaneously in one iteration:

$$\boldsymbol{\sigma}^{(t+1)} = \operatorname{sgn} \Big(\boldsymbol{\Xi} \, \boldsymbol{\Xi}^{\top} \boldsymbol{\sigma}^{(t)} \Big) \,.$$

Years later, Krotov & Hopfield (2016) improves the memory capacity of Hopfield network from $\mathcal{O}(d)$ to $\mathcal{O}(2^{d/2})$ when storing random samples. They adopted higher-order polynomial or exponential function (Demircigil et al., 2017) to distinguish each stored memory to alleviate the fuzzy memory problem: $\boldsymbol{\sigma}^{(t+1)} = \operatorname{sgn}(\boldsymbol{\Xi}(\boldsymbol{\Xi}^{\top}\boldsymbol{\sigma}^{(t)})^k)$ or $\boldsymbol{\sigma}^{(t+1)} = \operatorname{sgn}(\boldsymbol{\Xi}\exp(\boldsymbol{\Xi}^{\top}\boldsymbol{\sigma}^{(t)}))$.

This concept has evolved significantly and was extended to memories with continuous value. Modern Hopfield networks abstract retrieval as a one-iteration update (Ramsauer et al., 2021), and their retrieval dynamics $\mathcal{T}(\mathbf{x})$ can be unified under a three-step procedure (Millidge et al., 2022):

- (1) **Similarity** [$\mathbf{s} = \sin(\mathbf{\Xi}, \mathbf{x})$]: The query \mathbf{x} is compared against all stored patterns $\boldsymbol{\xi}_{1\cdots N}$ using the similarity function $\sin(\cdot, \cdot)$, obtaining a vector of similarity scores $\mathbf{s} \in \mathbb{R}^N$, where $\mathbf{s}_k = \sin(\boldsymbol{\xi}_k, \mathbf{x})$.
- (2) **Separation** [$\mathbf{p} = \operatorname{sep}(\mathbf{s})$]: The similarity scores (logits) \mathbf{s} are transformed by a separation function $\operatorname{sep}(\cdot)$ into a probability distribution \mathbf{p} . The separation function sharpens the scores and emphasizes patterns with high similarity.
- (3) **Readout** [$y = \Xi p$]: The final retrieved pattern y is computed as a weighted combination of stored memory patterns, using the weights p provided by the separation function.

Combining each step gives the unified retrieval dynamics:

$$\mathbf{y} = \mathcal{T}(\mathbf{x}) = \mathbf{\Xi} \operatorname{sep}(\operatorname{sim}(\mathbf{\Xi}, \mathbf{x})).$$

Concretely, Ramsauer et al. (2021) proposed using softmax(\cdot) function as the separation function, which further enlarges the memory capacity and draws a tight connection between associative memory and attention mechanism, with the retrieval dynamics being $\mathcal{T}(\mathbf{x}) = \Xi \operatorname{softmax}\left(\Xi^{\top}\mathbf{x}\right)$. Later, Hu et al. (2023) proposed sparse Hopfield network substituting $\operatorname{softmax}(\cdot)$ with $\operatorname{sparsemax}(\cdot)$, for inducing sparse selection while retaining differentiability. More recently, Wu et al. (2024a) attempted to store memory patterns in a kernel space with greater separation among patterns, giving rise to adding a new modulation step to the existing three-step unified framework. For clarity, we use the modulation function $\operatorname{mod}(\cdot)$ to describes how memory patterns are stored or pre-trained for better retrieval and larger capacity. So, it broadens the unified framework to:

$$\mathbf{y} = \mathbf{\Xi} \operatorname{sep}(\operatorname{sim}(\operatorname{mod}(\mathbf{\Xi}), \mathbf{x})).$$

The kernelized Hopfield network (Wu et al., 2024a) adopted $\operatorname{mod}(\Xi) = \Phi^\top \Phi \Xi$ for a learnable matrix $\Phi \in \mathbb{R}^{D_{\Phi} \times d}$, that projects memory patterns into a kernel space with the retrieval dynamics being $\mathcal{T}(\mathbf{x}) = \Xi \operatorname{sep}((\Phi \Xi)^\top (\Phi \mathbf{x})) = \Xi \operatorname{sep}((\Phi \Xi)^\top \mathbf{x})$. The kernel Φ is trained to minimize

a separation loss defined on Ξ , so that the expected Euclidean distance between any two memory patterns is maximized. A succeeding work (Hu et al., 2024) uses spherical codes to find the optimal kernel Ξ that maximizes the capacity of the kernelized Hopfield network.

Furthermore, the energy-based view is a defining feature of Hopfield networks: memory retrieval can be viewed as descending on a *energy landscape* (a Lyapunov function) $E(\cdot)$ whose minima coincide with stored patterns (or their modulated version). Formally, the retrieval dynamics $\mathcal{T}(\mathbf{x})$ and the corresponding energy function $E(\mathbf{x})$ are jointly and carefully designed such that each update monotonically decreases the energy (i.e., $E(\mathcal{T}(\mathbf{x})) < E(\mathbf{x})$), and successful retrieval occurs when being sufficiently close to a generalized fixed point near a specific memory pattern $\boldsymbol{\xi}_k \in \boldsymbol{\Xi}$ (i.e., $\|\mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}_k\|_2 < \epsilon$). This principled linkage between dynamics and energy ensures convergence and provides a powerful interpretable model of memory retrieving for Hopfield networks. Connecting to the previous unified framework, the separation function decides the direction of the retrieval dynamics $\mathcal{T}(\cdot)$, and the modulation function reshapes the geometry of the energy landscape $E(\cdot)$, and we organize all existing Hopfield networks' components and energy function in Table 1.

However, across existing formulations, the energy and retrieval dynamics are anchored to a fixed, task-agonistic similarity measure (typically the dot product). Apart from that, the energy and dynamics are solely determined by stored memories Ξ that overlook the nature of the subtle, nuanced, context-specific *association* between queries and memories required for "correct" retrieval. This fundamental gap motivates our work: to refine the energy and retrieval dynamics around a learnable, adaptive similarity measure while preserving the precious interpretability of Hopfield networks.

3 METHODS

In this section, we first establish a rigorous probabilistic framework to define *correct retrieval*, eliminating limitations of conventional proximity-based metrics (Section 3.1). To make this concept practical, we develop the similarity footprint (Section 3.2), a multi-dimensional descriptor measuring association between queries and memory patterns, and use it to learn an adaptive similarity integrated into an adaptive Hopfield network that achieves optimal correct retrieval for noisy, masked, and biased types of variants (Section 3.3).

3.1 VARIANT DISTRIBUTION AND CORRECT RETRIEVAL

Conventional analyses of associative memory (Ramsauer et al., 2021; Hu et al., 2023; Wu et al., 2024a; Hu et al., 2024) mostly focus on ϵ -retrieval:

```
Definition 1: \epsilon-retrieval (Hu et al., 2023; Wu et al., 2024a; Hu et al., 2025) Given a query \mathbf{x} \in \mathbb{R}^d and the retrieval result \mathbf{y} \in \mathbb{R}^d given by the memory system, a memory pattern \boldsymbol{\xi} \in \boldsymbol{\Xi} is said to be \epsilon-retrieved if \|\mathbf{y} - \boldsymbol{\xi}\|_2 \le \epsilon.
```

While ϵ -retrieval ensures the retrieval result \mathbf{y} lies near a certain stored memory pattern $\boldsymbol{\xi}_k \in \boldsymbol{\Xi}$, it provides no guarantee that $\boldsymbol{\xi}_k$ is the most appropriate match for query \mathbf{x} . The query \mathbf{x} may have stronger associations with a different pattern $\boldsymbol{\xi}_j$ ($j \neq k$), denoting that $\boldsymbol{\xi}_j$ could be the more appropriate match for \mathbf{x} . This identifies that proximity alone is an insufficient proxy for correctness.

To address this limitation, we use a probability distribution to model the generative process of the query \mathbf{x} . We posit that a query \mathbf{x} is not an arbitrary vector but a *variant* of a specific stored memory pattern $\boldsymbol{\xi} \in \boldsymbol{\Xi}$ generated by a context-dependent process. We formalize this via the *variant distribution*, which models the relation of memory patterns $\boldsymbol{\xi} \in \boldsymbol{\Xi}$ and queries $\mathbf{x} \in \mathbb{R}^d$ as variants:

Definition 2: Variant distribution

A variant distribution $\mathcal{V}(\Xi)$ is a joint distribution over pair $(\xi, \mathbf{x}) \in \Xi \times \mathbb{R}^d$ where $\xi \in \Xi$ is one of the stored memory patterns and $\mathbf{x} \in \mathbb{R}^d$ is an arbitrary query. For $(\xi, \mathbf{x}) \sim \mathcal{V}(\Xi)$, the probability density function $p_{\mathcal{V}(\Xi)}(\xi, \mathbf{x})$ (or $p_{\mathcal{V}}(\xi, \mathbf{x})$ when unambiguous) measures the likelihood of observing ξ and \mathbf{x} at the same time.

Additionally, the posterior $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$ represents the likelihood that query \mathbf{x} originates from memory pattern \mathbf{x} , and the likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ models how probable that $\boldsymbol{\xi}$ generates \mathbf{x} . This leads to a rigorous definition of the context-dependent *correct retrieval*:

Definition 3: Correct retrieval

A query x is said to be correctly retrieved under $\mathcal{V}(\Xi)$, if the retrieval result y satisfies that:

$$\underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,min}} \left\{ \|\mathbf{y} - \boldsymbol{\xi}'\|_{2} \right\} = \underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,max}} \left\{ p_{\mathcal{V}}(\boldsymbol{\xi}'|\mathbf{x}) \right\}. \tag{1}$$

In Equation 1, the left-hand side identifies the closest memory pattern to the retrieval result y (given by the memory system), while the right-hand side is ground truth (the most probable origin of query x given by variant distribution $\mathcal{V}(\Xi)$). Thus, intuitively, correct retrieval requires that the closest memory pattern coincides with ground truth. With further derivation,

$$\arg\max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \{ p_{\mathcal{V}}(\boldsymbol{\xi}'|\mathbf{x}) \} = \arg\max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}') \cdot \frac{p_{\mathcal{V}}(\boldsymbol{\xi}')}{p_{\mathcal{V}}(\mathbf{x})} \right\} = \arg\max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \{ p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}') \cdot p_{\mathcal{V}}(\boldsymbol{\xi}') \}. \quad (2)$$

This reformulation is necessary as modeling the likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ is more tractable than directly estimating the posterior $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$. The likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ is conditioned on a single, finite, known memory $\boldsymbol{\xi}$, while the posterior $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$ requires estimating a complex function that maps the entire query space \mathbb{R}^d to a discrete distribution over $\boldsymbol{\Xi}$. Given that the prior $p_{\mathcal{V}}(\boldsymbol{\xi})$ is typically uniform or can be easily estimated from samples, the central challenge of achieving correct retrieval reduces to accurately modeling $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$, in other words, how probable does \mathbf{x} generate $\boldsymbol{\xi}$ under $\mathcal{V}(\boldsymbol{\Xi})$? With this in hand, it is possible to model three canonical and common variant types rigorously:

Definition 4: Noisy variant

A query \mathbf{x} is a noisy variant if it is generated by adding Gaussian noise to a certain memory pattern $\boldsymbol{\xi} \in \boldsymbol{\Xi}$. Formally, $(\mathbf{x} - \boldsymbol{\xi}) \sim \mathcal{N}(\mathbf{0}, \operatorname{diag}(\boldsymbol{\sigma}))$ holds for $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}_{\operatorname{noisy}}(\boldsymbol{\Xi})$, where $\operatorname{diag}(\mathbf{v})$ transform vector bv to a diagonal matrix. The likelihood of noisy variant is:

$$p_{\mathcal{V}_{\text{noisy}}}(\mathbf{x}|\boldsymbol{\xi}) = \frac{1}{(2\pi)^{d/2}|\mathrm{diag}(\boldsymbol{\sigma})|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\xi})^{\top}\mathrm{diag}(\boldsymbol{\sigma})^{-1}(\mathbf{x} - \boldsymbol{\xi})\right).$$

Noisy variants have been widely studied (Krotov & Hopfield, 2016; Hu et al., 2023; Wu et al., 2024a), and it occurs in scenarios such as sensor noise. Specially, under isotropy $\sigma = \sigma \mathbf{1}$, the respective likelihood reduces to: $p_{\mathcal{V}_{noisy}}(\mathbf{x}|\boldsymbol{\xi}) = (2\pi\sigma)^{-d/2} \exp(-\|\mathbf{x} - \boldsymbol{\xi}\|_2^2/2\sigma)$.

Definition 5: Masked variant

A masked variant of a memory pattern $\xi \in \Xi$ is obtained by changing values in each dimension with probability p_{masked} to numbers generated by \mathcal{G} . The likelihood of masked variant is:

$$p_{\mathcal{V}_{\text{masked}}}(\mathbf{x}|\boldsymbol{\xi}) = \exp\left(\ln p_{\text{masked}} \cdot \sum_{i=1}^{d} [1 - \delta(\mathbf{x}_i - \boldsymbol{\xi}_i)]\right) \times \prod_{i=1}^{d} p_{\mathcal{G}}(\mathbf{x}_i)^{[1 - \delta(\mathbf{x}_i - \boldsymbol{\xi}_i)]}.$$

Masked variants arise in real-world scenarios such as information loss during transmission, the same object appearing in different background, and more.

Definition 6: Biased variant

 Adding a global bias to memory patterns gives the biased variant. Formally, $\mathbf{x} - \boldsymbol{\xi} = \mathbf{d}$ holds for $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}_{\text{biased}}(\boldsymbol{\Xi})$ and a constant vector $\mathbf{d} \in \mathbb{R}^d$. The likelihood of biased variant:

$$p_{\mathcal{V}_{\text{biased}}}(\mathbf{x}|\boldsymbol{\xi}) = \delta \left[d - \sum_{i=1}^d \delta(\mathbf{x}_i - \boldsymbol{\xi}_i - \mathbf{d}_i) \right], \quad \text{where } \delta(x) = \begin{cases} 1 & x = 0 \\ 0 & \text{otherwise} \end{cases}.$$

Biased variants occur as a systematic difference, such as changes in light conditions or use of filters.

We visualize the conditional probability density function $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ in Fig. 1, providing intuition akin to an *electron cloud*, with a memory pattern $\boldsymbol{\xi}$ as the atom nucleus and its variants as orbiting electrons. A direct observation is that $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ varies significantly across contexts, and may be analytically intractable. For instance, even though one can visualize the masked + noisy variant (Fig. 1d) by composing these two operations, deriving its likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ is analytically cumbersome. Consequently, although $\mathcal{V}(\boldsymbol{\Xi})$ is a principled tool to link queries with memory patterns, it poses two challenges for correct retrieval: (1) the underlying variant type is generally unknown a priori; and (2) the resulting variant distribution $\mathcal{V}(\boldsymbol{\Xi})$ can be too complex to model explicitly.

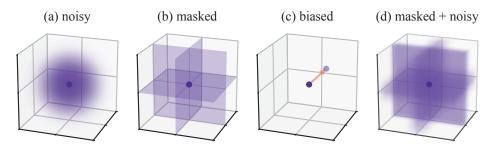


Figure 1: Visualization of probability density function $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ for noisy, masked, biased, and noisy + masked variants. Darker regions indicate larger $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$, and the central dark point represents $\boldsymbol{\xi}$.

3.2 SIMILARITY FOOTPRINT

In the last section, we established correct retrieval as selecting the pattern that maximizes $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$. This suggests that $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ itself serves as an ideal similarity function for associative memory, as it guides the memory system to return memory pattern chosen by $\arg\max_{\boldsymbol{\xi}\in\Xi}\{p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})\}$ (recall Equation 2), meeting the requirement of correct retrieval (Definition 3). Because $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ is often unknown and intractable, we instead mine for richer evidence from observable quantities to mimic the behavior of $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$. We introduce *similarity footprint*, which extracts a multi-dimensional descriptor of the relation between query \mathbf{x} and memory patterns. The core insight is that different generative processes leave distinct multi-dimensional signatures between a query and its origin memory pattern, suggesting replacing the scalar similarity measure with forming a convincing decision from the rich evidence described by the structured descriptor.

We begin with a *base similarity* (e.g., dot product $\boldsymbol{\xi}^{\top} \mathbf{x}$ or Euclidean distance $-\|\boldsymbol{\xi} - \mathbf{x}\|_2$), and define the *k-optimal similarity* between $\boldsymbol{\xi}$ and \mathbf{x} :

$$\sin^{(k)}(\boldsymbol{\xi},\mathbf{x}) \triangleq \max_{D \subseteq [d], |D| = k} \left\{ \sin(\boldsymbol{\xi}_D,\mathbf{x}_D) \right\}, \quad \text{where } \mathbf{v}_D \triangleq \left[\mathbf{v}_{D_1}, \mathbf{v}_{D_2}, \cdots, \mathbf{v}_{D_{|D|}} \right]^\top.$$

Here, \mathbf{v}_D is a sub-vector of \mathbf{v} containing only the elements corresponding to indices in D. Intuitively, $\sin_k(\boldsymbol{\xi}, \mathbf{x})$ quantifies the best agreement between \mathbf{x} and $\boldsymbol{\xi}$ focusing only on their k most similar dimensions. This allows the similarity measure to ignore potentially corrupted dimensions and delve into the most informative dimensions. Base on this, we then define the *similarity footprint* as the vector of these k-optimal similarity $\sin^{(k)}(\cdot,\cdot)$ across all possible dimensionalities:

$$\operatorname{ftpt}_{\operatorname{sim}}(\boldsymbol{\xi}, \mathbf{x}) \triangleq \left[\sin^{(d)}(\boldsymbol{\xi}, \mathbf{x}), \ \sin^{(d-1)}(\boldsymbol{\xi}, \mathbf{x}), \ \cdots, \ \sin^{(1)}(\boldsymbol{\xi}, \mathbf{x}) \right]^{\top}$$

This serves as the rich descriptor of the relation between \mathbf{x} and $\boldsymbol{\xi}$, providing more evidence for measuring similarity. However, a naïve computation of the footprint requires evaluating all 2^d-1 subspaces, which is impractical. Fortunately, for decomposable similarity functions (such as dot product and Euclidean distance), whose results can be computed by aggregating similarity in each dimension, the footprint $\operatorname{ftpt}_{\operatorname{sim}}(\cdot,\cdot)$ can be obtained efficiently in $\mathcal{O}(d\log d)$ time by sorting. Let \mathbf{q} be the dimension-wise similarity vector, where $\mathbf{q}_i = \sin(\boldsymbol{\xi}_i, \mathbf{x}_i)$ for $i \in [d]$, and let $\tilde{\mathbf{q}}$ be the vector \mathbf{q} sorted in ascending order. Then, the similarity footprint calculation is equivalent to:

$$ftpt_{sim}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{U}\tilde{\mathbf{q}}.$$
 (3)

where **U** is the upper-right triangle matrix of $\mathbf{1}_{d\times d}$, (i.e., $\mathbf{U}_{i,j}=1$ if $1\leq i\leq j\leq d$, and $\mathbf{U}_{i,j}=0$ otherwise). This operation literally calculates the cumulative sum of the sorted dimension-wise similarities, as finding the k-optimal similarity is equivalent to aggregating the k largest dimension-wise similarities for decomposable base similarity.

3.3 ADAPTIVE SIMILARITY AND ADAPTIVE HOPFIELD NETWORK

The similarity footprint provides a structured, multi-scale descriptor of the association between a query \mathbf{x} and a memory pattern $\boldsymbol{\xi}$. To leverage these crucial evidences and create a similarity measure that adapts to the underlying variant distribution, we define a learnable *adaptive similarity* as a linear function of the footprint: $s_{\text{sim}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \text{ftpt}_{\text{sim}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} \tilde{\mathbf{q}}$, for some learnable weight vector $\mathbf{w} \in \mathbb{R}^d$ and a base similarity $\sin(\cdot, \cdot)$. This formulation allows the model to learn the

relative importance of similarity across different subspaces and focus on the informative ones. For instance, for masked variants, the model might assign higher weights to the last d-m terms in the footprint, incorporating the uncorrupted information, whereas for noisy variants, it might distribute weights to larger subspaces for a global view.

To further enhance the model's expressiveness, we can combine footprints of multiple base similarities and derive the final similarity function $s(\xi, \mathbf{x})$ and its vectorized form $\mathbf{s}(\xi, \mathbf{x})$:

$$s(\boldsymbol{\xi}, \mathbf{x}) = \sum_{k=1}^{B} \beta_k \cdot \mathbf{w}_k^{\top} \operatorname{ftpt}_{\operatorname{sim}_k}(\boldsymbol{\xi}, \mathbf{x}) \quad \text{ and } \quad \mathbf{s}(\boldsymbol{\Xi}, \mathbf{x}) = \left[s(\boldsymbol{\xi}_1, \mathbf{x}), \ s(\boldsymbol{\xi}_2, \mathbf{x}), \ \cdots, \ s(\boldsymbol{\xi}_N, \mathbf{x}) \right]^{\top},$$

for some learnable scalars $\beta_{1...B}$, and B different base similarities. In this work, we use two simple and common measures as base similarities: the Euclidean distance $\operatorname{dis}(\boldsymbol{\xi},\mathbf{x}) = -\|\boldsymbol{\xi} - \mathbf{x}\|_2^2$ and dot product $\operatorname{dot}(\boldsymbol{\xi},\mathbf{x}) = \boldsymbol{\xi}^{\top}\mathbf{x}$. This adaptive similarity can be seamlessly integrated into the modern unified Hopfield network framework (Ramsauer et al., 2021; Millidge et al., 2022) by incorporating a separation function $\operatorname{sep}(\cdot,\cdot)$, and we name the resulting model as *adaptive Hopfield network* (A-Hop). When using $\operatorname{softmax}(\cdot)$ as the separation function, the retrieval dynamics of A-Hop is:

$$\mathbf{y} = \mathcal{T}(\mathbf{x}) = \mathbf{\Xi} \operatorname{sep}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x})) = \mathbf{\Xi} \operatorname{softmax}(\beta_1 \cdot \mathbf{s}_{\operatorname{dis}}(\mathbf{\Xi}, \mathbf{x}) + \beta_2 \cdot \mathbf{s}_{\operatorname{dot}}(\mathbf{\Xi}, \mathbf{x}))$$
 (4)

The parameters \mathbf{w} 's and β 's are optimized using a sample set drawn from the variant distribution. The learning objective is to minimize the discrepancy between the model's predicted likelihood $\tilde{p}_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}) \triangleq \sup(\mathbf{s}(\boldsymbol{\Xi},\mathbf{x}))$ and the underlying ground-truth likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$. However, since $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_i)$ is unknown, we use $\delta(\boldsymbol{\xi}_i - \boldsymbol{\xi})$ for $i \in [N]$ (a one-hot vector) as an approximation, and minimize the cross-entropy loss:

$$\mathcal{L}(\Xi, \mathcal{V}) = \mathbb{E}_{(\mathbf{x}, \boldsymbol{\xi}) \sim \mathcal{V}(\Xi)} \left[-\sum_{i=1}^{N} \delta(\boldsymbol{\xi}_{i} - \boldsymbol{\xi}) \log \tilde{p}_{\mathcal{V}}(\mathbf{x} | \boldsymbol{\xi}) \right]$$
 (5)

Furthermore, A-Hop achieves optimal correct retrieval (Definition 7) for noisy, masked, and biased variants when weights w's are decided ideally. Additionally, A-Hop guarantees a decreasing and bounded energy when the retrieval process is iterative (more than one step). These are formalized in the following theorems, with detailed proofs provided in Appendix A.2.

Definition 7: Optimal correct retrieval

We say a retrieval dynamics $\mathcal{T}(\mathbf{x})$ achieves optimal correct retrieval under $\mathcal{V}(\Xi)$, if for any $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\Xi)$ it achieves correct retrieval for query \mathbf{x} .

Theorem 1: A-Hop retrieval dynamics

The following retrieval dynamics adopted by A-Hop achieves optimal correct retrieval for noisy, masked, and biased variants, with a careful design of $s(\Xi, x)$:

$$\mathbf{y} = \mathcal{T}(\mathbf{x}) = \mathbf{\Xi} \operatorname{sep}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x}))$$

Theorem 2: A-Hop energy landscape

Energy $E(\mathbf{x})$ will be monotonically decreasing and its value could be bounded for isotropic noisy, and biased variants, if the following energy is used:

$$E(\mathbf{x}) = -\mathrm{lse}\left(\mathbf{s}(\mathbf{\Xi}, \mathbf{x})\right)$$

4 EXPERIMENTS

We assess the effectiveness of A-Hop on tasks including memory retrieval, tabular classification, image classification, and multiple instance learning, demonstrating that A-Hop achieves state-of-the-art performance on these tasks. A further ablation study validates our design choice of similarity footprint. Due to space constraints, full descriptions of baselines, metrics, datasets, and implementation details are provided in Appendix A.4.

4.1 Memory Retrieval

Prior work on Hopfield networks primarily assesses retrieval accuracy (Definition 9) under two settings: (1) masking half of the dimensions (masked variant); and (2) adding Gaussian noise (noisy

variant). To more comprehensively probe retrieval robustness, we introduce mixed variants parameterized by a triplet $(d_{\text{mask}}, d_{\text{noise}}, d_{\text{bias}}) \in [0, 1]^3$ that controls the intensity (difficulty) of masking, noise, and bias, respectively (see Appendix A.4.3 for formal definitions).

We evaluate A-Hop and baselines on 12 mixed-variant types (Fig. 2c) with 64-dimensional random memory patterns at scales of 2048 (Fig. 2a) and 4096 (Fig. 2b). We further assess high intensity of masked, noisy, and biased combined settings on 2048 synthetic vectors and on 2048 samples from MNIST as memory patterns (Table 2). While baseline models perform well in noisy settings, their accuracy degrades sharply when faced with more complex variant settings. In contrast, A-Hop's adaptive similarity enables it to maintain high retrieval accuracy and low retrieval error across all tested scenarios. This highlights its robustness and impressive adaptability to align similarity to the underlying variant distribution through learning.

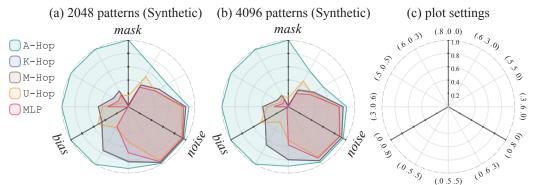


Figure 2: Retrieval accuracy (†) for different variant settings

Table 2: Retrieval accuracy (\uparrow) and error (\downarrow) between models. Each cell contains the mean accuracy or error with standard deviation in a smaller font. Results of the best-performing model are **bolded**. Difficulty d means that variant setting ($d_{\text{mask}}, d_{\text{noise}}, d_{\text{bias}}$) = (d, d, d) is used.

Dataset	Synthetic				MNIST			
Difficulty	0.4 0.5		5	0.6		0.7		
Metrics	Accuracy	Error	Accuracy	Error	Accuracy	Error	Accuracy	Error
М-Нор	.520±.02	.176±.01	.195±.03	.300±.01	.875±.01	.013±.00	.661±.02	.068±.00
U-Hop	.260±.04	$.417 {\scriptstyle \pm .02}$	$.059 \pm .01$	$.554 {\pm} .01$.540±.03	$.143 \pm .01$.176±.02	$.347 \pm .01$
K-Hop	.487±.03	$.295 {\scriptstyle \pm .02}$	$.177 \pm .02$	$.764 {\scriptstyle \pm .02}$.764±.02	$.064 \pm .01$.526±.02	$.164 {\scriptstyle \pm .01}$
K_2 -Hop 1	.521±.02	$.176 \pm .01$	$.195 \pm .02$	$.298 \pm .01$.878±.01	$.013 \pm .00$.660±.01	$.068 \pm .01$
А-Нор	.724 ±.02	.106 ±.01	.360 ±.02	.227 ±.01	.939 ±.01	$\textbf{.005} {\pm}.00$.849 ±.01	.015 ±.01

4.2 TABULAR CLASSIFICATION

Table 3: Predictive performance (↑) between models on tabular data. Each cell contains mean accuracy or AUC-ROC score with standard deviation in a smaller font. Results of the best-performing associative memory are **bolded**, and the best other model is <u>underlined</u>.

Model	Adult	Bank	Vaccine	Purchase	Heart
М-Нор К ₂ -Нор А-Нор	.8080±.001 .8172±.003 .8634±.002	.9085±.003 .9092±.002 .9139±.002	.7975±.001 .7971±.003 .8042 ±.002	.8822±.001 .8825±.002 .9007 ±.001	.6325±.002 .6473±.002 .7315±.002
Extra Trees Random Forest AdaBoost XGBoost	$.8595 \pm .004 \\ .8592 \pm .002 \\ .8597 \pm .003 \\ \underline{.8640} \pm .002$	$.9098 \pm .003 \\ .9132 \pm .003 \\ .9094 \pm .001 \\ \underline{.9152} \pm .003$	$.7932 \pm .002 \\ .7918 \pm .003 \\ .8011 \pm .002 \\ \underline{.8034} \pm .002$	$.8916 \pm .002 \\ .9002 \pm .001 \\ .8865 \pm .001 \\ \underline{.9032} \pm .003$	$.7175 \pm .003$ $.7254 \pm .002$ $.7294 \pm .001$ $\underline{.7370} \pm .003$

 $^{^{1}}$ K₂-Hop is K-Hop whose kernel is optmized by Equation 5, rather than the original separation loss.

We integrate Hopfield networks into a memory-based classifier (see Appendix A.4.4), and evaluated predictive performance of A-Hop and baselines on five tabular datasets. Table 3 shows that A-Hop consistently outperforms all other associative memory-based classifiers, and demonstrates an advantage over Extra Trees (Geurts et al., 2006), Random Forest (Breiman, 2001), and AdaBoost (Freund & Schapire, 1997). However, XGBoost (Chen & Guestrin, 2016) has a slight edge over A-Hop on 4 out of 5 datasets, dominating this task.

Notably, Adult and Heart appear harder than the other datasets: A-Hop yields a 5-10% absolute gain over memory-based baselines on these two, compared with less than 2% on others. This pattern suggests that these datasets exhibit subtler, heterogeneous variant distributions (e.g., mixed feature types, sparsity, complexity) where adaptive similarity better aligns the memory's neighborhood structure with the data geometry. This validates the potential of adaptivity in this domain.

4.3 IMAGE CLASSIFICATION AND MULTIPLE INSTANCE LEARNING

Table 4: Classification accuracy (†) of each model on images, and AUC-ROC score (†) of each model in multiple instance learning task. Each cell contains accuracy or AUC-ROC score with standard deviation in a smaller font. Results of the best-performing associative memory are **bolded**.

Image Classification				Multipl	le Instance L	earning		
Dataset	CIFAR10	CIFAR100	Tiny ImageNet	Dataset	Tiger	Fox	Elephant	UCSB
М-Нор	.5123±.003	.2464±.003	.1095±.002	М-Нор	.8924±.005	.6327±.013	.9344±.009	.8815±.022
K-Hop	.5489±.002	$.2877 \pm .002$	$.1164 {\pm}.002$	S-Hop	.8923±.006	$.6433 \pm .015$	$.9365 \pm .002$	$.8794 \pm .024$
A-Hop	.5637±.003	.2904 ±.002	$\textbf{.}1213 \pm .002$	A-Hop	.9030 ±.007	.6753 ±.013	.9451 ±.004	$\pmb{.8935} {\pm .022}$

Following established protocols, we evaluate A-Hop on image classification (Wu et al., 2024a) and multiple instance learning (Ramsauer et al., 2021; Hu et al., 2023) by integrating it as a layer within larger and more complicated neural network architectures (i.e., HopfieldLayer, HopfieldPooling). As shown in Table 4, A-Hop consistently achieves the highest scores among all Hopfield variants in both task categories. This demonstrates that the benefits of adaptive similarity extend to complex, high-dimensional data and can enhance the performance of sophisticated models like the Image Transformer. While the performance gains are more modest compared to the memory retrieval task, this is expected, as the HopfieldLayer is one component within a much larger model. Nevertheless, the consistent improvement confirms that optimizing the similarity measure remains a valuable factor for enhancing performance in complex deep learning systems.

4.4 ABLATION STUDY

Due to page limit, the ablation study is moved to Appendix A.4.6.

5 Conclusion

We reframe associative memory retrieval as a problem of *correct retrieval* under a task- and context-dependent variant distribution, motivating a similarity measure that approximates the likelihood that a stored pattern generated the query. Building on this principle, we propose adaptive similarity, prove its optimality for three canonical variant families (noisy, masked, biased), and instantiate it in a new adaptive Hopfield network, A-Hop. This perspective clarifies why fixed, pre-defined similarities are inherently limited: they cannot align to the prevailing variant distribution and thus struggle to guarantee correctness, whereas adaptivity enables the model to capture the underlying variant distribution through samples, shifting towards correctness.

Empirically, A-Hop establishes state-of-the-art performance among Hopfield networks across memory retrieval, tabular classification, image classification, and multiple instance learning. The gains are most pronounced under mixed variant settings where adaptive similarity maintains impressively high retrieval accuracy and low error. In downstream tasks, A-Hop consistently improves over prior Hopfield variants. Ultimately, adaptive similarity is a key principle for advancing associative memories, paving the way for more powerful and resilient memory systems.

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. In this study, no human subjects or animal experimentation was involved. All datasets used were sourced in compliance with relevant usage guidelines, ensuring no violation of privacy. We have taken care to avoid any biases or discriminatory outcomes in our research process. No personally identifiable information was used, and no experiments were conducted that could raise privacy or security concerns. We are committed to maintaining transparency and integrity throughout the research process.

7 REPRODUCIBILITY STATEMENT

Code We provide code to help understand this work, and is publicly available at: https://anonymous.4open.science/r/Adaptive-Hopfield-Network-C137/.

Datasets All datasets are either included in the repo, or a description for how to download and preprocess the dataset is provided. All datasets are public and raise no ethical concerns.

Hyperparameters All parameters of our proposed framework are in Appendix A.4.

Environment Details of our experimental setups are provided in Appendix A.4.

Random Seed we do not set a random seed specifically for all random behavior, with the random seed determined PyTorch.

8 LLM USAGE

Large Language Models (LLMs) were used to aid polishing of the manuscript. Specifically, we used an LLM to assist in refining the language, improving readability, and ensuring clarity in various sections of the paper. The model helped with tasks such as sentence rephrasing, grammar checking, and enhancing the overall flow of the text.

It is important to note that the LLM was not involved in the ideation, research methodology, or experimental design. All research concepts, ideas, and analyses were developed and conducted by the authors. The contributions of the LLM were solely focused on improving the linguistic quality of the paper, with no involvement in the scientific content or data analysis.

REFERENCES

Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: https://doi.org/10.24432/C5XW20.

Leo Breiman. Random Forests. *Machine Learning*, 45(1):5–32, October 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL https://link.springer.com/article/10.1023/A:1010933404324. Company: Springer Distributor: Springer Institution: Springer Label: Springer Number: 1 Publisher: Kluwer Academic Publishers.

P. Bull, I. Slavitt, and G. Lipstein. Harnessing the power of the crowd to increase capacity for data science in the social sector. In *ICML #Data4Good Workshop*, 2016.

Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings* of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, pp. 785–794, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939785. URL https://doi.org/10.1145/2939672.2939785.

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(2):288–299, July 2017. ISSN 1572-9613. doi: 10.1007/s10955-017-1806-y. URL https://doi.org/10.1007/s10955-017-1806-y.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

```
DrivenData. Flu shot learning: Predict h1n1 and seasonal flu vaccines. https://www.drivendata.org/competitions/66/flu-shot-learning/data/, 2019.
```

- Yoav Freund and Robert E. Schapire. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997. ISSN 0022-0000. doi: https://doi.org/10.1006/jcss.1997.1504. URL https://www.sciencedirect.com/science/article/pii/S002200009791504X.
- Elisa Drelie Gelasca, Jiyun Byun, Boguslaw Obara, and B.S. Manjunath. Evaluation and benchmark for biological image segmentation. In *IEEE International Conference on Image Processing*, Oct 2008. URL http://vision.ece.ucsb.edu/publications/elisa_ICIP08.pdf.
- Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*, 63(1):3–42, April 2006. ISSN 1573-0565. doi: 10.1007/s10994-006-6226-1. URL https://doi.org/10.1007/s10994-006-6226-1.
- J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, April 1982. doi: 10.1073/pnas.79.8.2554. URL https://www.pnas.org/doi/abs/10.1073/pnas.79.8.2554. Publisher: Proceedings of the National Academy of Sciences.
- Jerry Yao-Chieh Hu, Donglin Yang, Dennis Wu, Chenwei Xu, Bo-Yu Chen, and Han Liu. On Sparse Modern Hopfield Model. *Advances in Neural Information Processing Systems*, 36:27594–27608, December 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/hash/57bc0a850255e2041341bf74c7e2b9fa-Abstract-Conference.html.
- Jerry Yao-Chieh Hu, Dennis Wu, and Han Liu. Provably Optimal Memory Capacity for Modern Hopfield Models: Transformer-Compatible Dense Associative Memories as Spherical Codes. *Advances in Neural Information Processing Systems*, 37:70693–70729, December 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/82846e19e6d42ebfd4ace4361def29ae-Abstract-Conference.html.
- Jerry Yao-Chieh Hu, Bo-Yu Chen, Dennis Wu, Feng Ruan, and Han Liu. Nonparametric modern hopfield models. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=xkV3uCQtJm.
- Kaggle. Cardiovascular disease dataset. https://www.kaggle.com/datasets/sulianova/cardiovascular-disease-dataset, 2018.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Dmitry Krotov and John J. Hopfield. Dense Associative Memory for Pattern Recognition. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://papers.nips.cc/paper_files/paper/2016/hash/eaae339c4d89fc102edd9dbdb6a28915-Abstract.html.
- Dmitry Krotov, Benjamin Hoover, Parikshit Ram, and Bao Pham. Modern Methods in Associative Memory, July 2025. URL http://arxiv.org/abs/2507.06211.arXiv:2507.06211 [cs].
- Gert Lanckriet and Bharath K. Sriperumbudur. On the Convergence of the Concave-Convex Procedure. In *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009. URL https://papers.nips.cc/paper_files/paper/2009/hash/8b5040a8a5baf3e0e67386c2e3a9b903-Abstract.html.
- Yann Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
- Beren Millidge, Tommaso Salvatori, Yuhang Song, Thomas Lukasiewicz, and Rafal Bogacz. Universal Hopfield Networks: A General Framework for Single-Shot Associative Memory Models. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 15561–15583. PMLR, June 2022. URL https://proceedings.mlr.press/v162/millidge22a.html. ISSN: 2640-3498.

- Yasushi Miyashita. Neuronal correlate of visual associative long-term memory in the primate temporal cortex. *Nature*, 335(6193):817–820, October 1988. ISSN 1476-4687. doi: 10.1038/335817a0. URL https://www.nature.com/articles/335817a0. Publisher: Nature Publishing Group.
- Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems*, 62:22–31, 2014. ISSN 0167-9236. doi: https://doi.org/10.1016/j.dss.2014.03.001. URL https://www.sciencedirect.com/science/article/pii/S016792361400061X.
- John M. Pearce and Mark E. Bouton. Theories of Associative Learning in Animals. *Annual Review of Psychology*, 52(Volume 52, 2001):111–139, February 2001. ISSN 0066-4308, 1545-2085. doi: 10.1146/annurev.psych.52.1.111. URL https://www.annualreviews.org/content/journals/10.1146/annurev.psych.52.1.111. Publisher: Annual Reviews.
- Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas Gruber, Markus Holzleitner, Thomas Adler, David Kreil, Michael K Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. Hopfield networks is all you need. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=tL89RnzIiCd.
- Bishwajit Saha, Dmitry Krotov, Mohammed J Zaki, and Parikshit Ram. End-to-end differentiable clustering with associative memories. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 29649–29670. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/saha23a.html.
- C. Sakar and Yomi Kastro. Online Shoppers Purchasing Intention Dataset. UCI Machine Learning Repository, 2018. DOI: https://doi.org/10.24432/C5F88Q.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
- Jane X. Wang, Lynn M. Rogers, Evan Z. Gross, Anthony J. Ryals, Mehmet E. Dokucu, Kelly L. Brandstatt, Molly S. Hermiller, and Joel L. Voss. Targeted enhancement of cortical-hippocampal brain networks and associative memory. *Science*, 345(6200):1054–1057, August 2014. doi: 10.1126/science.1252900. URL https://www.science.org/doi/full/10.1126/science.1252900. Publisher: American Association for the Advancement of Science.
- Shurong Wang, Zhuoyang Shen, Xinbao Qiao, Tongning Zhang, and Meng Zhang. Dynfrs: An efficient framework for machine unlearning in random forest. In Y. Yue, A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *International Conference on Representation Learning*, volume 2025, pp. 10636–10657, 2025. URL https://proceedings.iclr.cc/paper_files/paper/2025/hash/1caf09c9f4e6b0150b06a07e77f2710c-Abstract-Conference.html.
- Dennis Wu, Jerry Yao-Chieh Hu, Teng-Yun Hsiao, and Han Liu. Uniform Memory Retrieval with Larger Capacity for Modern Hopfield Models. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 53471–53514. PMLR, July 2024a. URL https://proceedings.mlr.press/v235/wu24i.html. ISSN: 2640-3498.
- Yu-Hsuan Wu, Jerry Yao-Chieh Hu, Weijian Li, Bo-Yu Chen, and Han Liu. STan-Hop: Sparse Tandem Hopfield Model for Memory- Enhanced Time Series Prediction. *International Conference on Representation Learning*, 2024:30886-30925, May 2024b. URL https://proceedings.iclr.cc/paper_files/paper/2024/hash/832b20b65f655587e9c0447860406a82-Abstract-Conference.html.

A APPENDIX

APPENDIX CONTENTS

A.1	Notatio	ons	14
A.2	Theore	ems	15
	A.2.1	Optimal Correct Retrieval	15
	A.2.2	Unified Adaptive Similarity and Energy Function	19
A.3	Discus	sion	21
	A.3.1	On Optimal Correct Retrieval	21
	A.3.2	On More Complicated Variant Distribution	22
A.4	Experi	ments	22
	A.4.1	Baselines and Metrics	22
	A.4.2	Datasets	23
	A.4.3	Memory Retrieval	24
	A.4.4	Tabular Classification	24
	A.4.5	Image Classification and Multiple Instance Learning	24
	A.4.6	Ablation Study	25

A.1 NOTATIONS

Table 5: Notations and symbolds used in this work.

	Table 5: Notations and symbolds used in this work.
Symbol	Description
$oldsymbol{\xi},oldsymbol{\xi}_k$ $oldsymbol{\Xi}$ d N	A specific memory pattern $(d \times 1)$ (or memory, stored memory pattern, stored pattern). The $d \times N$ memory matrix, with each memory pattern being its column vector. The dimensionality of memory patterns. The number of stored memory patterns.
$sim(\boldsymbol{\xi}, \mathbf{x})$	The similarity function that measures how strong the association are between the inputs (or similarity, similarity measure, measure, association).
$sep(\mathbf{s}) \\ mod(\mathbf{\Xi}) \\ E(\mathbf{x}) \\ \mathcal{T}(\mathbf{x}) \\ \mathbf{p} \\ \mathbf{s}$	The separation function, turning the output of $\operatorname{sim}(\cdot,\cdot)$ (logits) to a probability distribution. The modulation function that governs how memory patterns are stored and learned. The energy landscape, defined on the same vector space as memory patterns. The retrieval dynamics, defined on the same vector space as memory patterns. The probability distribution vector produced by $\operatorname{sep}(\cdot)$. The similarity score vector produced by $\operatorname{sim}(\cdot,\cdot)$.
x	The query vector. Also, the input to the associative memory
$oldsymbol{y}{\mathcal{V}(oldsymbol{\Xi})}$	The retrieval result vector. Also, the output of the associative memory. The variant distribution on memory matrix Ξ , governs how queries are generated. Each query \mathbf{x} is sampled from this distribution together with its origin memory pattern $\boldsymbol{\xi}$. The joint probability density function that measures the likelihood that $\boldsymbol{\xi}$ and \mathbf{x} are
$p_{\mathcal{V}}(\boldsymbol{\xi}, \mathbf{x})$ $p_{\mathcal{V}}(\boldsymbol{\xi} \mathbf{x})$	observed together. The conditional probability density function (posterior) that measures the likelihood that
$p_{\mathcal{V}}(\mathbf{x} \mathbf{x})$ $p_{\mathcal{V}}(\mathbf{x} \mathbf{\xi})$	originates from \mathbf{x} when observed \mathbf{x} . The joint probability density function (likelihood) that measures the likelihood that $\boldsymbol{\xi}$ generates \mathbf{x} when observed $\boldsymbol{\xi}$.
q	The dimension-wise similarity vector whose value of the <i>i</i> -th index measures the similarity between the value of <i>i</i> -th index in \mathbf{x} and $\boldsymbol{\xi}$.
$egin{array}{c} ilde{\mathbf{q}} \ extbf{U} \end{array}$	The sorted version of q (sorted in ascending order). The upper right triangle matrix of ones.
$\mathrm{dis}(oldsymbol{\xi},\mathbf{x})$	The (negative and squared) Euclidean distance similarity $-\ \mathbf{x} - \boldsymbol{\xi}\ _2^2$.
$dot(\boldsymbol{\xi}, \mathbf{x})$ $sim^{(k)}(\boldsymbol{\xi}, \mathbf{x})$	The dot product similarity $\mathbf{x}^{\top}\boldsymbol{\xi}$. The k -optimal similarity function that finds a k -dimensional subspace that maximizes the similarity $\sin(\cdot, \cdot)$ of the inputs within that subspace.
$\mathrm{ftpt}_{\mathrm{sim}}(oldsymbol{\xi},\mathbf{x})$	The similarity footprint function that generates the rich descriptor between ξ and x with $sim(\cdot, \cdot)$ being the base similarity (or footprint).
$s_{\mathrm{sim}}(oldsymbol{\xi},\mathbf{x})$	The adaptive similarity function adopting $\operatorname{ftpt}_{\operatorname{sim}}(\cdot,\cdot)$ with $\operatorname{sim}(\cdot,\cdot)$ as the base similarity
$\mathbf{s}_{\mathrm{sim}}(\mathbf{\Xi},\mathbf{x})$	The vectorized form of the adaptive similarity function $s_{\text{sim}}(\cdot, \cdot)$, and returns a vector that measures the adaptive similarity between $\boldsymbol{\xi}_i$ and \mathbf{x} for $i \in [N]$.
$s({m \xi},{f x})$	The final adaptive similarity function that aggregate multiple $s_{\text{sim}_k}(\cdot, \cdot)$ for different base similarity $\text{sim}(\cdot, \cdot)$ / footprint.
$\mathbf{s}(\mathbf{\Xi},\mathbf{x})$	The vectorized form of the final adaptive similarity function $s(\cdot, \cdot)$, and returns a vector that measures the adaptive similarity between ξ_i and \mathbf{x} for $i \in [N]$.
w	The weight vector that turns the footprint into a scalar, which is designed to extract information from the rich descriptor.
$egin{aligned} eta \ \mathcal{L}(oldsymbol{\xi}, \mathcal{V}) \end{aligned}$	Scalar used to aggregate different adaptive similarities $\mathbf{s}_{\text{sim}}(\cdot, \cdot)$. The loss function used to optimize \mathbf{w} 's and β 's
[n]	The set of integers less than or equal to n .
$\operatorname{sgn}(x)$	Return the sign $(-1 \text{ or } +1)$ of the input.
$egin{array}{c} \delta(x) \ \mathbf{v}^{ op} \end{array}$	The Dirac delta that returns 1 when the input is 0 and returns 0 otherwise.
	Transpose of a vector / matrix.
diam'(TT)	Transform vector \mathbf{v} to a diagonal matrix.
$\operatorname{diag}(\mathbf{v})$	A sub waster of recontaining only the elements companding to indices in D
$\mathbf{v}_D = \ \mathbf{v}\ _p$	A sub-vector of \mathbf{v} containing only the elements corresponding to indices in D . The ℓ_p norm.

A.2 THEOREMS

We define the retrieval accuracy that estimates the retrieval performance of an associative memory under a certain variant distribution.

Definition 8: Retrieval accuracy

Retrieval accuracy for an associative memory with retrieval dynamics $\mathcal{T}(\cdot)$ is the probability that correct retrieval is met:

$$\mathbb{E}_{(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})} \left[\delta \left(\underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg min}} \{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_{2} \} - \underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg max}} \{ p_{\mathcal{V}}(\boldsymbol{\xi}' | \mathbf{x}) \} \right) \right]$$

$$= \Pr_{(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})} \left[\underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg min}} \{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_{2} \} = \underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg max}} \{ p_{\mathcal{V}}(\boldsymbol{\xi}' | \mathbf{x}) \} \right]$$

However, Definition 8 is usually intractable as $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ is unknown and complicated. Therefore, we define empirical retrieval accuracy based on samples drawn from $\mathcal{V}(\boldsymbol{\Xi})$, which is computable, and used in our experiments (Section 4).

Definition 9: Empirical retrieval accuracy

Empirical retrieval accuracy for an associative memory with retrieval dynamics $\mathcal{T}(\cdot)$ can be estimated by performing abundant retrieval tests:

$$\begin{split} & \mathbb{E}_{(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})} \left[\delta \bigg(\underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg\,min}} \{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_2 \} - \boldsymbol{\xi} \bigg) \right] \\ & = \Pr_{(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})} \left[\underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\operatorname{arg\,min}} \{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_2 \} = \boldsymbol{\xi} \right] \end{split}$$

A.2.1 OPTIMAL CORRECT RETRIEVAL

We now start to prove Theorem 1.

Let us begin with a simple variant — the isotropic noisy variant.

Lemma 1: Optimal correct retrieval for isotropic noisy variant

A-Hop achieves optimal correct retrieval (Definition 7) for isotrophic noisy variant $\mathcal{V}_{noisy}(\Xi)$ (Definition 4, $(\mathbf{x} - \boldsymbol{\xi}) \sim \mathcal{N}(\mathbf{0}, \sigma \mathbf{I})$ for $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}_{noisy}(\Xi)$ some $\sigma \in \mathbb{R}$) for arbitrary memory matrix $\Xi \in \mathbb{R}^{d \times N}$.

Proof. We claim that the optimal correct retrieval is achieved when using $sep(\cdot) = arg max(\cdot)$, and $ftpt_{dis}(\xi, \mathbf{x})$ only (i.e., $\beta_1 = 1$ and $\beta_2 = 0$ in Equation 4). That is, the retrieval dynamics should be:

$$\begin{split} \mathcal{T}(\mathbf{x}) &= \underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\arg\max} \left\{ \mathbf{w}^{\top} \mathrm{ftpt}_{\mathrm{dis}}(\boldsymbol{\xi}', \mathbf{x}) \right\} \\ &= \underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\arg\max} \left\{ -\|\boldsymbol{\xi}' - \mathbf{x}\|_{2}^{2} \right\} \end{split}$$

This step can be satisfied by setting $\mathbf{w}_1 = 1$, and $\mathbf{w}_i = 0$ for $2 \le i \le d$. Then, optimal retrieval is achieved only when Equation 1 is met. We first estimate the right-hand side of Equation 1:

$$\begin{split} \arg\max_{\boldsymbol{\xi}'\in\boldsymbol{\Xi}} \left\{p_{\mathcal{V}_{\text{noisy}}}(\boldsymbol{\xi}|\mathbf{x})\right\} &= \argmax_{\boldsymbol{\xi}'\in\boldsymbol{\Xi}} \left\{\ln p_{\mathcal{V}_{\text{noisy}}}(\mathbf{x}|\boldsymbol{\xi})\right\} \\ &= \arg\max_{\boldsymbol{\xi}'\in\boldsymbol{\Xi}} \left\{-\frac{d}{2}\ln 2\pi\sigma - \frac{1}{2\sigma}\|\mathbf{x} - \boldsymbol{\xi}'\|_2^2\right\} \\ &= \arg\max_{\boldsymbol{\xi}'\in\boldsymbol{\Xi}} \left\{-\|\boldsymbol{\xi}' - \mathbf{x}\|_2^2\right\} \end{split}$$

The first step comes from Equation 2, and assuming that the prior $p(\xi)$ is uniform (which is often the case for memory retrieval) or can be easily obtained from samples. And the second step comes

from Definition 4. We can see that the derived results coincide with the retrieval dynamics derived before. Therefore, plugging the retrieval dynamics to the left-hand side of Equation 1 gives:

$$\underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,min}} \left\{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_{2} \right\} = \underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,min}} \left\{ \left\| \underset{\boldsymbol{\xi}''}{\operatorname{arg\,max}} \left\{ -\| \boldsymbol{\xi}'' - \mathbf{x} \|_{2}^{2} \right\} - \boldsymbol{\xi}' \right\|_{2} \right\}$$

$$= \sum_{k=1}^{N} \delta \left(\left\| \underset{\boldsymbol{\xi}''}{\operatorname{arg\,max}} \left\{ -\| \boldsymbol{\xi}'' - \mathbf{x} \|_{2}^{2} \right\} - \boldsymbol{\xi}_{k} \right\|_{2} \right) \cdot \boldsymbol{\xi}_{k}$$

$$= \sum_{k=1}^{N} \delta \left(\underset{\boldsymbol{\xi}''}{\operatorname{max}} \left\{ -\| \boldsymbol{\xi}'' - \mathbf{x} \|_{2}^{2} \right\} - \left[-\| \boldsymbol{\xi}_{k} - \mathbf{x} \|_{2}^{2} \right] \right) \cdot \boldsymbol{\xi}_{k}$$

$$= \underset{\boldsymbol{\xi}_{k} \in \Xi}{\operatorname{arg\,max}} \left\{ -\| \boldsymbol{\xi}_{k} - \mathbf{x} \|_{2}^{2} \right\}$$

The second step holds as there always exists a $\xi' \in \Xi$ that let $\|\arg\max_{\xi''} \{-\|\xi'' - \mathbf{x}\|_2^2\} - \xi'\|_2 = 0$, since the resulting vector of the $\arg\max(\cdot) \in \Xi$, and ξ' iterates every column vector of Ξ , thus must have coincided with resulting vector, and the thrid step holds for a similar reason.

Therefore, we show that the left-hand side and right-hand side of Equation 1 are the same $(\arg\max_{\boldsymbol{\xi}'\in\Xi}\{-\|\boldsymbol{\xi}_k-\mathbf{x}\|_2^2\})$. Thus, the requirement for correct retrieval is met for all $(\boldsymbol{\xi},\mathbf{x})\sim\mathcal{V}(\Xi)$, yielding optimal correct retrieval.

If we adopt a footprint that does not sort the dimension-wise similarity vector ${\bf q}$ by substituting $\tilde{{\bf q}}$ in Equation 3 to ${\bf q}$ and gives $\operatorname{ftpt}_{\operatorname{dis}'}({\boldsymbol \xi},{\bf x})={\bf U}{\bf q}$, we can prove the optimality for the standard noisy variant defined in Definition 4, which is more general than Lemma 1. However, the footprint $\operatorname{ftpt}_{\operatorname{dis}}({\boldsymbol \xi},{\bf x})={\bf U}\tilde{{\bf q}}$ achieves high empirical retrieval accuracy, but it is harder to estimate analytically.

Lemma 2: Optimal retrieval for noisy variant

A-Hop achieves optimal correct retrieval (Definition 7) for noisy variant $\mathcal{V}_{\text{noisy}}(\Xi)$ (Definition 4 for arbitrary memory matrix $\Xi \in \mathbb{R}^{d \times N}$.

Proof. Following the same spirit in the proof of Lemma 1. One can see that the right-hand side (RHS) of Equation 1 is (similar to Lemma 1):

$$\underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,max}} \left\{ p_{\mathcal{V}_{\text{noise}}}(\boldsymbol{\xi}'|\mathbf{x}) \right\} = \underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,max}} \left\{ -\frac{1}{2} \ln \left[(2\pi)^d | \operatorname{diag}(\boldsymbol{\sigma})| \right] - \frac{1}{2} (\mathbf{x} - \boldsymbol{\xi}')^\top \operatorname{diag}(\boldsymbol{\sigma})^{-1} (\mathbf{x} - \boldsymbol{\xi}') \right\} \\
= \underset{\boldsymbol{\xi}' \in \Xi}{\operatorname{arg\,max}} \left\{ -\sum_{i=1}^d \frac{(\boldsymbol{\xi}_i' - \mathbf{x})^2}{\boldsymbol{\sigma}_i} \right\}$$

While the left-hand side (LHS) of Equation 1 is (step 1 follows how RHS is resolved in Lemma 1):

$$\underset{\boldsymbol{\xi}' \in \Xi}{\arg \min} \left\{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_{2} \right\} = \underset{\boldsymbol{\xi}_{k} \in \Xi}{\arg \max} \left\{ \mathbf{w}^{\top} \mathbf{U} \left(\boldsymbol{\xi}_{k} - \mathbf{x} \right)^{2} \right\} \\
= \underset{\boldsymbol{\xi}_{k} \in \Xi}{\arg \max} \left\{ \mathbf{u}^{\top} (\boldsymbol{\xi}_{k} - \mathbf{x})^{2} \right\} \\
= \underset{\boldsymbol{\xi}_{k} \in \Xi}{\arg \max} \left\{ \sum_{i=1}^{d} \mathbf{u}_{i} (\boldsymbol{\xi}_{k,i} - \mathbf{x}_{i})^{2} \right\}$$

Here, $\mathbf{v}^2 \in \mathbb{R}^d$ denotes a dimension-wise square operation over $\mathbf{v} \in \mathbb{R}^d$, so that $(\boldsymbol{\xi}_k - \mathbf{x})_i^2 = (\boldsymbol{\xi}_{k,i} - \mathbf{x}_i)^2$. Also, we let $\mathbf{u}^\top = \mathbf{w}^\top \mathbf{U}$ for simplicity. One can see that LHS equals RHS when:

$$\forall i, i \in [d], \mathbf{u}_i = -\frac{1}{\sigma_i}$$

Since U is a full-rank matrix, so that $\mathbf{w}^{\top} = \mathbf{u}^{\top} \mathbf{U}^{-1}$ holds. When we set w in the following way:

$$\mathbf{w}_i = \begin{cases} -\boldsymbol{\sigma}_1^{-1} & i = 1\\ \boldsymbol{\sigma}_{i-1}^{-1} - \boldsymbol{\sigma}_i^{-1} & 2 \leq i \leq d \end{cases}$$

LHS and RHS of Equation 1 are the same, satisfying the requirement for correct retrieval. Furthermore, we can tell that Lemma 1 is a special case of this lemma.

Lemma 3: Optimal retrieval for masked variant

 A-Hop achieves optimal correct retrieval (Definition 7) for masked variant $\mathcal{V}_{\text{masked}}(\Xi)$ (Definition 5 for arbitrary memory matrix $\Xi \in \mathbb{R}^{d \times N}$, and for a uniform generator \mathcal{G} ($p_{\mathcal{G}}(\cdot)$ is a constant).

Proof. As in Lemma 1, we first reformulate the RHS of Equation 1, and find the suitable choice for w to make LHS of Equation 1 equal RHS. We let \mathbf{q} be the dimension-wise similarity vector with $\mathbf{q}_i = -(\boldsymbol{\xi}_i - \mathbf{x}_i)^2$, and $\tilde{\mathbf{q}}$ the vector that sort \mathbf{q} in ascending order. The RHS can be expanded as:

$$\begin{aligned} & \underset{\boldsymbol{\xi}' \in \Xi}{\arg\max} \left\{ p_{\mathcal{V}_{\text{masked}}}(\boldsymbol{\xi}'|\mathbf{x}) \right\} \\ & = \underset{\boldsymbol{\xi}' \in \Xi}{\arg\max} \left\{ \ln p_{\text{masked}} \cdot \sum_{i=1}^{d} [1 - \delta(\mathbf{x}_i - \boldsymbol{\xi}_i')] + \sum_{i=1}^{d} [1 - \delta(\mathbf{x}_i - \boldsymbol{\xi}_i)] \ln p_{\mathcal{G}}(\mathbf{x}_i') \right\} \\ & = \underset{\boldsymbol{\xi}' \in \Xi}{\arg\max} \left\{ \left(\ln p_{\text{masked}} + \ln p_{\mathcal{G}} \right) \cdot \sum_{i=1}^{d} [1 - \delta(\mathbf{x}_i - \boldsymbol{\xi}_i')] \right\} \\ & = \underset{\boldsymbol{\xi}' \in \Xi}{\arg\max} \left\{ -d + \sum_{i=1}^{d} \delta(\mathbf{x}_i - \boldsymbol{\xi}_i') \right\} \\ & = \underset{\boldsymbol{\xi}' \in \Xi}{\arg\max} \left\{ \sum_{i=1}^{d} \delta(\tilde{\mathbf{q}}_i) \right\} \end{aligned}$$

Here, the second step is valid as $p_{\mathcal{G}}$ is a constant, and we term this constant $p_{\mathcal{G}}$. As $p_{\mathcal{G}}$ is a constant, it must be less than or equal to 1 to make $p_{\mathcal{G}}(\cdot)$ a valid probability density function. Additionally, we know $0 \le p_{\text{masked}} < 1$ from definition, and therefore, $\ln p_{\text{masked}} + \ln p_{\mathcal{G}} < 0$, and this explains why a sign change occurs in step three. The derived RHS suggests designing a discrete adaptive similarity:

$$s_{\mathrm{dis}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\delta}(\tilde{\mathbf{q}}_i) = \sum_{i=1}^{d} \mathbf{w}_i \delta(\tilde{\mathbf{q}}_i)$$

Then, the LHS would be (first step following that of Lemma 1, and recall we use $\arg \max(\cdot)$ as separation function):

$$\underset{\boldsymbol{\xi}' \in \boldsymbol{\Xi}}{\arg\min} \left\{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_2 \right\} = \underset{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}}{\arg\max} \left\{ \mathbf{w}^\top \boldsymbol{\delta}(\tilde{\mathbf{q}}_i) \right\} = \underset{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}}{\arg\max} \left\{ \sum_{i=1}^d \mathbf{w}_i \cdot \delta(\tilde{\mathbf{q}}_i) \right\}$$

Setting $\mathbf{w} = \mathbf{1}$ concludes that LHS equals RHS for all $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})$, and thus, the optimal correct retrieval is achieved for this concrete adaptive similarity $s_{\mathrm{dis}}(\boldsymbol{\xi}, \mathbf{x})$.

It can be shown that it is impossible to find a continuous $s_{\rm dis}(\boldsymbol{\xi},\mathbf{x})$ for masked variant's optimal correct retrieval, unless more constraints on $\boldsymbol{\xi}$ and \mathbf{x} are made. Typically, such a continuous function is possible if $\|\boldsymbol{\xi} - \mathbf{x}\|_2 \ge \varepsilon$ (can be bounded from below) for $\varepsilon > 0$.

Lemma 4: Optimal retrieval for biased variant

A-Hop achieves optimal correct retrieval (Definition 7) for biased variant $\mathcal{V}_{\text{biased}}(\Xi)$ (Definition 5 for arbitrary memory matrix $\Xi \in \mathbb{R}^{d \times N}$, and an arbitrary difference vector $\mathbf{d} \in \mathbb{R}^d$.

Proof. Following the proof to Lemma 1, the LHS of Equation 1 is:

$$\arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ p_{\mathcal{V}_{\text{biased}}}(\boldsymbol{\xi}' | \mathbf{x}) \right\} = \arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ \delta \left[d - \sum_{i=1}^{d} \delta(\mathbf{x}_{i} - \boldsymbol{\xi}'_{i} - \mathbf{d}_{i}) \right] \right\}$$

$$= \arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ -d + \sum_{i=1}^{d} \delta(\mathbf{x}_{i} - \boldsymbol{\xi}'_{i} - \mathbf{d}_{i}) \right\}$$

$$= \arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ -\|\mathbf{x} - \boldsymbol{\xi}' - \mathbf{d}\|_{2}^{2} \right\}$$

The last step follows that the maximum score is both 0 before and after the transform, and the goal is to assign a high score (0) when $\mathbf{x} - \boldsymbol{\xi} = \mathbf{d}$. Here, we use a similar continuous adaptive similarity defined in the main text:

$$s_{\mathrm{dis}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} \mathbf{q} - \mathbf{q}^{\top} \mathbf{q}$$

with $\mathbf{q} = \mathbf{x} - \boldsymbol{\xi}$, and set $\mathbf{u}^{\top} = \mathbf{w}^{\top} \mathbf{U}$ with $\mathbf{u} = 2\mathbf{d}^{\top}$. Then, the RHS of Equation 1 is:

$$\begin{aligned} \operatorname*{arg\,min}_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \left\{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_2 \right\} &= \operatorname*{arg\,max}_{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}} \left\{ \mathbf{u}^\top \mathbf{q} - \mathbf{q}^\top \mathbf{q} \right\} \\ &= \operatorname*{arg\,max}_{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}} \left\{ 2\mathbf{d}^\top \mathbf{q} - \mathbf{q}^\top \mathbf{q} - \mathbf{d}^\top \mathbf{d} \right\} \\ &= \operatorname*{arg\,max}_{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}} \left\{ -(\mathbf{q} - \mathbf{d})^\top (\mathbf{q} - \mathbf{d}) \right\} \\ &= \operatorname*{arg\,max}_{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}} \left\{ -\|\mathbf{q} - \mathbf{d}\|_2^2 \right\} \\ &= \operatorname*{arg\,max}_{\boldsymbol{\xi}_k \in \boldsymbol{\Xi}} \left\{ -\|\mathbf{x} - \boldsymbol{\xi}_k - \mathbf{d}\|_2^2 \right\} \end{aligned}$$

This follows immediately that LHS equals RHS, and the optimal correct retrieval is achieved when:

$$\mathbf{w}_i = \begin{cases} 2\mathbf{d}_1 & i = 1 \\ 2\mathbf{d}_i - 2\mathbf{d}_{i-1} & 2 \le i \le d \end{cases}$$

It finally comes down to Theorem 1.

Theorem 1: A-Hop retrieval dynamics

The following retrieval dynamics adopted by A-Hop achieves optimal correct retrieval for noisy, masked, and biased variants, with a careful design of $s(\Xi, x)$:

$$\mathbf{y} = \mathcal{T}(\mathbf{x}) = \mathbf{\Xi} \operatorname{sep}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x}))$$

Proof. First of all, $sep(\cdot) = arg max(\cdot)$ is crucial for achieving optimal correct retrieval, as it transforms the left-hand side of Equation 1 as (see Lemma 1):

$$\underset{\boldsymbol{\xi}' \in \Xi}{\arg \min} \left\{ \| \mathcal{T}(\mathbf{x}) - \boldsymbol{\xi}' \|_2 \right\} = \underset{\boldsymbol{\xi}_k \in \Xi}{\arg \max} \left\{ \mathbf{s}_{\mathrm{dis}}(\boldsymbol{\xi}_k, \mathbf{x}) \right\}$$

In Lemma 2, we see that using the following adaptive similarity achieves optimal correct retrieval:

$$s_{\text{dis}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} \mathbf{q}$$
 with $\mathbf{w}_i = \begin{cases} -\boldsymbol{\sigma}_1^{-1} & i = 1\\ \boldsymbol{\sigma}_{i-1}^{-1} - \boldsymbol{\sigma}_i^{-1} & 2 \le i \le d \end{cases}$

In Lemma 3, we see that using the following adaptive similarity achieves optimal correct retrieval:

$$s_{\mathrm{dis}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \boldsymbol{\delta}(\tilde{\mathbf{q}})$$
 with $\mathbf{w}_i = 1$ for $i \in [d]$

In Lemma 4, we see that using the following adaptive similarity achieves optimal correct retrieval:

$$s_{\mathrm{dis}}(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} (\mathbf{x} - \boldsymbol{\xi}) - (\mathbf{x} - \boldsymbol{\xi})^{\top} (\mathbf{x} - \boldsymbol{\xi}) \qquad \text{with } \mathbf{w}_i = \begin{cases} 2\mathbf{d}_1 & i = 1\\ 2\mathbf{d}_i - 2\mathbf{d}_{i-1} & 2 \leq i \leq d \end{cases}$$

One can see that achieving optimal correct retrieval is not easy, and it requires the sacrifice of continuous. However, we can build a continuous adaptive similarity inspired from the proof of Theorem 1 that achieve high retrieval accuracy (at least, empirically). For more discussion on this topic, please read Appendix A.3.1.

A.2.2 Unified Adaptive Similarity and Energy Function

We can find that the adaptive similarity in Lemma 1 has the form:

$$s(\boldsymbol{\xi}, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi})^{\top} (\mathbf{x} - \boldsymbol{\xi})$$

while that of Lemma 2 has the form:

$$s(\boldsymbol{\xi}, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi})^{\top} \operatorname{diag}(\mathbf{a})(\mathbf{x} - \boldsymbol{\xi})$$

for some diagonal matrix $\operatorname{diag}(\mathbf{a})$, and $\mathbf{a}_i > 0$ for all $i \in [d]$. Meanwhile, for Lemma 4 has the form:

$$s(\boldsymbol{\xi}, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi})^{\top} (\mathbf{x} - \boldsymbol{\xi}) + \mathbf{b}^{\top} (\mathbf{x} - \boldsymbol{\xi})$$

for some real vector b. That is being said that we can unifies these three adaptive similarity by:

$$s_{\text{unify}}(\boldsymbol{\xi}, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi})^{\top} \operatorname{diag}(\mathbf{a})(\mathbf{x} - \boldsymbol{\xi}) + \mathbf{b}^{\top}(\mathbf{x} - \boldsymbol{\xi})$$

However, this similarity is too tough, we can analysis a simpler one:

$$s_{\text{unify}}(\boldsymbol{\xi}, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi})^{\top}(\mathbf{x} - \boldsymbol{\xi}) + \mathbf{b}^{\top}(\mathbf{x} - \boldsymbol{\xi})$$

If we use a $\operatorname{softmax}(\cdot)$ function as the separation function and $s_{\operatorname{unify}}(\boldsymbol{\xi}, \mathbf{x})$ as the similarity function, and construct an energy function (with $s_{\operatorname{unify}}(\boldsymbol{\Xi}, \mathbf{x})$ being the vectorized form of $s_{\operatorname{unify}}(\boldsymbol{\xi}, \mathbf{x})$):

$$E(\mathbf{x}) = -\mathrm{lse}(\mathbf{s}_{\mathrm{unify}}(\mathbf{\Xi}, \mathbf{x})) \tag{6}$$

whose gradient is:

$$\nabla_{\mathbf{x}} E(\mathbf{x}) = -\operatorname{softmax}(\mathbf{s}_{\text{unify}}(\mathbf{\Xi}, \mathbf{x}))^{\top} \nabla_{\mathbf{x}} \mathbf{s}_{\text{unify}}$$

Letting $\mathbf{p}_i(\mathbf{x}) \triangleq \operatorname{softmax}(\mathbf{s}_{\text{unify}}(\mathbf{\Xi}, \mathbf{x}))_i$:

$$\nabla_{\mathbf{x}} E(\mathbf{x}) = -\sum_{i=1}^{N} \mathbf{p}_{i}(\mathbf{x}) \nabla_{\mathbf{x}} s_{\text{unify}}(\boldsymbol{\xi}_{i}, \mathbf{x})$$

$$= -\sum_{i=1}^{N} \mathbf{p}_{i}(\mathbf{x}) \cdot (-2\mathbf{x} + 2\boldsymbol{\xi}_{i} + \mathbf{b})$$

$$= 2\mathbf{x} - \mathbf{b} - 2\sum_{i=1}^{N} \mathbf{p}_{i}(\mathbf{x}) \cdot \boldsymbol{\xi}_{i}$$

Retrieval on the gradient flow gives:

$$\frac{d\mathbf{x}}{dt} = -\nabla_{\mathbf{x}} E(\mathbf{x}) = -2\mathbf{x} + \mathbf{b} + 2\sum_{i=1}^{N} \mathbf{p}_{i}(\mathbf{x}) \cdot \boldsymbol{\xi}_{i}$$

Then, consider using gradient descent with step η , where $\eta > 0$ for discrete-time retrieval:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla_{\mathbf{x}} E(\mathbf{x}^{(t)})$$
$$= (1 - 2\eta) \cdot \mathbf{x}^{(t)} + \eta \mathbf{b} + 2\eta \sum_{i=1}^{N} \mathbf{p}_{i}(\mathbf{x}^{(t)}) \cdot \boldsymbol{\xi}_{i}$$

By setting $\eta = \frac{1}{2}$ that would cancels the x term on the RHS and remove the coefficient before the summation, which is wonderful:

$$\mathbf{x}^{(t+1)} = \frac{1}{2}\mathbf{b} + \mathbf{\Xi}\operatorname{softmax}(\mathbf{s}_{\text{unify}}(\mathbf{\Xi}, \mathbf{x}))$$

From Lemma 4, we know that the setting $\mathbf{b} = 2\mathbf{d}$ is optimal for noisy variant, plugging this in gives:

$$\mathbf{x}^{(t+1)} - \mathbf{d} = \mathcal{T}(\mathbf{x}^{(t)}) = \mathbf{\Xi} \operatorname{softmax}(\mathbf{s}_{\text{unify}}(\mathbf{\Xi}, \mathbf{x}))$$
 (7)

suggesting that we add a new de-bias term $-\mathbf{d}$ for biased variants, which coincidentally, remove the bias vector \mathbf{d} . However, when there is no bias, Equation 7 reduce to the simple retrieval dynamics we are familiar with.

We then further analysis the behavior of the energy (Equation 6) retrieval dynamics Equation 7:

Lemma 5: Rewriting the energy

Let the energy function be $E(\mathbf{x}) = -\mathrm{lse}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x}))$, where $s(\boldsymbol{\xi}_i, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi}_i)^{\top}(\mathbf{x} - \boldsymbol{\xi}_i) + \mathbf{b}^{\top}(\mathbf{x} - \boldsymbol{\xi}_i)$. Then, the energy can be written as

$$E(\mathbf{x}) = \|\mathbf{x}\|_2^2 - \operatorname{lse}(\mathbf{A}\mathbf{x} + \mathbf{c})$$

for some matrix $\mathbf{A} \in \mathbb{R}^{N \times d}$ and vector $\mathbf{c} \in \mathbb{R}^{N}$.

Proof.

$$E(\mathbf{x}) = -\operatorname{lse}\left(\mathbf{s}(\mathbf{\Xi}, \mathbf{x})\right)$$

$$= -\operatorname{ln}\sum_{i=1}^{N} \exp\left[-\|\mathbf{x}\|_{2}^{2} - \|\boldsymbol{\xi}_{i}\|_{2}^{2} + (2\boldsymbol{\xi}_{i} + \mathbf{b})^{\top}\mathbf{x} + \boldsymbol{\xi}_{i}^{\top}\mathbf{b}\right]$$

$$= -\operatorname{ln}\sum_{i=1}^{n} \exp(-\|\mathbf{x}\|_{2}^{2}) + \exp\left[(2\boldsymbol{\xi}_{i} + \mathbf{b})^{\top}\mathbf{x} + \boldsymbol{\xi}_{i}^{\top}\mathbf{b} - \|\boldsymbol{\xi}_{i}\|_{2}^{2}\right]$$

$$= \operatorname{ln}n + \|\mathbf{x}\|_{2}^{2} - \operatorname{lse}(\mathbf{A}\mathbf{x} + \mathbf{c})$$

for $\mathbf{A}_i^{\top} = 2\boldsymbol{\xi}_i + \mathbf{b}$ and $\mathbf{c}_i = \boldsymbol{\xi}_i^{\top} \mathbf{b} - \|\boldsymbol{\xi}_i\|_2^2$. Also, we can omit the term $\ln n$ as it is a constant. \square

Therefore, we have decomposed the energy function $E(\mathbf{x})$ into a convex function $g(\mathbf{x}) = ||\mathbf{x}||_2^2$, and a concave function $-h(\mathbf{x}) = -\operatorname{lse}(\mathbf{A}\mathbf{x} + \mathbf{c})$.

Lemma 6: Decreasing energy

Energy function $E(\mathbf{x})$ would be monotonically decreasing using the retrieval dynamics:

$$\mathbf{x}^{(t+1)} - \mathbf{d} = \mathcal{T}(\mathbf{x}^{(t)}) = \mathbf{\Xi} \operatorname{softmax}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x})) \quad \text{for } s(\boldsymbol{\xi}_i, \mathbf{x}) = -(\mathbf{x} - \boldsymbol{\xi}_i)^{\top}(\mathbf{x} - \boldsymbol{\xi}_i) + \mathbf{b}^{\top}(\mathbf{x} - \boldsymbol{\xi}_i)$$

Proof. Using the concave convex procedure (Lanckriet & Sriperumbudur, 2009), we construct a convex surrogate function $U_t(\mathbf{x})$ for each iteration t by linearizing the concave function $-h(\mathbf{x})$ around the current $\mathbf{x}^{(t)}$:

$$U_t(\mathbf{x}) = g(x) - \left[h(\mathbf{x}^{(t)}) + \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)})^{\top} (\mathbf{x} - \mathbf{x}^{(t)}) \right]$$

= $\|\mathbf{x}\|_2^2 - \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)})^{\top} \mathbf{x} + (\nabla_{\mathbf{x}} h(\mathbf{x}^{(t)})^{\top} \mathbf{x}^{(t)} - h(\mathbf{x}^{(t)}))$

The next $\mathbf{x}^{(t+1)}$ is the minimizer of $U_t(\mathbf{x})$, i.e., $\mathbf{x}^{(t+1)} = \arg\min_{\mathbf{x} \in \mathbb{R}^d} \{U_t(\mathbf{x})\}$, and we can find it by setting its gradient to zero:

$$\nabla_{\mathbf{x}} U_t(\mathbf{x}) = 2\mathbf{x} - \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)}) = \mathbf{0} \implies \mathbf{x}^{(t+1)} = \frac{1}{2} \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)}) = \frac{1}{2} \mathbf{A}^{\top} \operatorname{softmax}(\mathbf{A} \mathbf{x}^{(t)} + \mathbf{c})$$

We can add the term $-\|\mathbf{x}\|_2^2$ back to $\operatorname{softmax}(\cdot)$ as it is independent of index i. Thus, by denoting $\mathbf{p}_i(\mathbf{x}_i) = \operatorname{softmax}(\mathbf{A}\mathbf{x}^{(t)} + \mathbf{c}) = \operatorname{softmax}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x}))$ (follows Lemma 5), we have:

$$\mathbf{x}^{(t+1)} = \frac{1}{2} \sum_{i=1}^{N} \mathbf{A}_i \cdot \mathbf{p}_i(\mathbf{x}^{(t)})$$

$$= \sum_{i=1}^{N} \boldsymbol{\xi}_i \cdot \mathbf{p}_i(\mathbf{x}^{(t)}) - \frac{1}{2} \mathbf{b} \sum_{i=1}^{N} \mathbf{p}_i(\mathbf{x}^{(t)})$$

$$= \mathbf{\Xi} \operatorname{softmax}(\mathbf{s}(\mathbf{\Xi}, \mathbf{x}^{(t)})) - \frac{1}{2} \mathbf{b}$$

and this agrees with what we have derived before (Equation 7).

Then, by convexity of $h(\mathbf{x})$ (recall $-h(\mathbf{x})$ is concave), we have the following inequality:

$$h(\mathbf{x}) \ge h(\mathbf{x}^{(t)}) + \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)})^{\top} (\mathbf{x} - \mathbf{x}^{(t)})$$

$$\implies g(\mathbf{x}) - h(\mathbf{x}) \le g(\mathbf{x}) - h(\mathbf{x}^{(t)}) + \nabla_{\mathbf{x}} h(\mathbf{x}^{(t)})^{\top} (\mathbf{x} - \mathbf{x}^{(t)}) = U_t(\mathbf{x})$$

with the equality holds iff $\mathbf{x} = \mathbf{x}^{(t)}$, and recall that $\mathbf{x}^{(t+1)}$ is the minimum value of $U_t(\mathbf{x})$, we have:

$$E(\mathbf{x}^{(t+1)}) \le U_t(\mathbf{x}^{(t+1)}) \le U_t(\mathbf{x}^{(t)}) = E(\mathbf{x}^{(t)})$$
(8)

with equality holds when $\mathbf{x}^{(t)} = \mathbf{x}^{(t+1)}$. Thus, $E(\mathbf{x}^{(t+1)}) \leq E(\mathbf{x}^{(t)})$ completes the proof.

We can see that $E(\mathbf{x}) = \|\mathbf{x}\|_2^2 - \mathcal{O}(\|\mathbf{x}\|_2)$, so that $E(\mathbf{x}) \to +\infty$ as $\|\mathbf{x}\|_2 \to +\infty$, so $E(\mathbf{x})$ is coercive, meaning its level sets are compact. Additionally, $U_t(\mathbf{x})$ is 2-strongly convex as $\nabla_{\mathbf{x}}^2 U_t(\mathbf{x}) = 2\mathbf{I}$, therefore,

$$U_t(\mathbf{x}^{(t)}) - U_t(\mathbf{x}^{t+1}) \ge \|\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}\|_2^2$$

Along with Inequality 8:

$$E(\mathbf{x}^{(t)}) - E(\mathbf{x}^{(t+1)}) \ge U_t(\mathbf{x}^{(t)}) - U_t(\mathbf{x}^{(t+1)}) \ge ||\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}||_2^2$$

This yields that

$$E(\mathbf{x}^{(t)}) - E(\mathbf{x}^{(0)}) \ge \sum_{\tau=0}^{t-1} \|\mathbf{x}^{(\tau+1)} - \mathbf{x}^{(\tau)}\|_2^2$$

This mean that $E(\mathbf{x}^{(t)})$ is bounded from below. We can see that the sequence of \mathbf{x} : $\{\mathbf{x}^{(t)}\}$ must remain within the level set $\{\mathbf{x} \mid E(\mathbf{x}) \leq E(\mathbf{x}^{(0)})\}$ as the sequence $E(\mathbf{x}^{(0)})$ is non-increasing. Therefore, $\{\mathbf{x}^{(t)}\}$ must remains in the compact set $\{\mathbf{x} \mid E(\mathbf{x}) \leq E(\mathbf{x}^{(0)})\}$, and be bounded.

Theorem 2: A-Hop energy landscape

energy Energy $E(\mathbf{x})$ will be monotonically decreasing and its value could be bounded for isotropic noisy, and biased variants, if the following energy is used:

$$E(\mathbf{x}) = -\mathrm{lse}\left(\mathbf{s}(\mathbf{\Xi}, \mathbf{x})\right)$$

Proof. In Lemma 6, we have proven that the energy function will be monotonically decreasing.

Also, from above analysis we can see that $E(\mathbf{x}^{(t)})$ could be bounded by:

$$0 \le \sum_{\tau=0}^{t-1} \|\mathbf{x}^{(\tau+1)} - \mathbf{x}^{(\tau)}\|_2^2 \le E(\mathbf{x}^{(t)}) \le E(\mathbf{x}^{(0)})$$

We try to find the property for a very different energy landscape and a different retrieval dynamics, also these retrieval dynamics can guaratee optimal correct retrieval, for iostropic noisy, and biased variant, when the separation function is $\arg\max(\cdot)$ and has weight **b** choosen ideally. Theorem 2 tries to connect the noval correct retrieval and the traditional energy analysis.

A.3 DISCUSSION

A.3.1 ON OPTIMAL CORRECT RETRIEVAL

Achieving optimal correct retrieval is costly, as it requires designing "weird" adaptive similarity function or adopting discrete function making the model unlearnable. For instance, it is impossible to achieve optimal correct retrieval for masked variants, as there would always be lower-bound issue for continuous functions. Also, using $\operatorname{softmax}(\cdot)$ as the separation make it hard to prove whether optimal retrieval is achieved or not, as we cannot exclude the effect of other memory patherns from the one receive largest similarity score.

However, we propose choosing the value of \mathbf{w} wisely for adaptive similarity $s(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} \tilde{\mathbf{q}}$ to achieve "sub-optimal" correct retrieval (well, it is hard to rigorously define sub-optimality).

For noisy variant, a choice that setting $\mathbf{w}_1 = 1$ and other values to 0 is suggested. This enable the model to have global view (mauniplate similarity in the largest subspace), and its electron cloud would look like a shpere as illustrated in Fig. 1a.

Next, for masked variant, we suggest using $\mathbf{w}_i = iC$ for some large constant $C > 2d^2$, because this punishes patterns that has very limited dimension where $\boldsymbol{\xi}_i = \mathbf{x}_i$, and stress importance on small subspaces. By doing so the similarity function will look like the likelihood electron cloud in Fig. 1b.

Finally, for biased variant, the best way to set \mathbf{w} is to set $\mathbf{u}^{\top} = \mathbf{w}^{\top}\mathbf{U}$ to $\tilde{\mathbf{d}}$, the sorted bias vector, so that they similarity function can focus on how comparing the difference of $\mathbf{x} - \boldsymbol{\xi}$ with a unkown but almost known bias.

The main point of this section is that the optimality of correct retrieval is too strict so that achieving so force us to abandon good properties. Also, in most cases, some very corner cases set very high difficulty for achieve optimal. However, correct retrieval itself is a good property, but making every retrieval correct is too strict. Therefore, an open question is that how to define a sub-optimal correct retrieval standard so that it guarantees great memory retrieval performance and leave us freedom.

A.3.2 ON MORE COMPLICATED VARIANT DISTRIBUTION

Variant distributions discussed so far are actually simple. In previous analysis, we assumes that all memory variants follows a similar distribution. For example, $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}_{\text{noisy}}(\boldsymbol{\Xi})$ ensures that all memory patterns generates a noisy variants, sharing a similar $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$. We say a variant distribution general if knowing $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_i)$ is equivalent to knowing $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_j)$ for arbitrary $i, j \in [N]$ and $i \neq j$. In other words, we can obtain $p_{\mathcal{V}}(\mathbf{x}|\mathbf{x}_j)$ by substituting \mathbf{x}_j with \mathbf{x}_i . However, there could be cases where each memory patterns generates variants quite differently. Intuitively, we say each pattern is generates on their own, meaning that $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi})$ has completely different form for different memory patterns, and we call such memory variant isolated.

By looking closely to the adaptive similarity function:

$$s(\boldsymbol{\xi}, \mathbf{x}) = \mathbf{w}^{\top} \mathbf{U} \tilde{\mathbf{q}}$$

One can see that it uses a universal weight \mathbf{w} for all pairs of $(\boldsymbol{\xi}, \mathbf{x}) \sim \mathcal{V}(\boldsymbol{\Xi})$, assuming that the variant distribution it is trying to model is general. However, such limitation can be easily broken by introducing more weights. For instance, we use separate weights for different memory patterns, i.e., for N memory patterns, we spare weight \mathbf{w}_k to memory pattern $\boldsymbol{\xi}_k$. Therefore, the adaptive similarity that can suit isolated variant distribution look like:

$$s(\boldsymbol{\xi}_k, \mathbf{x}) = \mathbf{w}_k^{\mathsf{T}} \mathbf{U} \tilde{\mathbf{q}}$$

This can fit each isolated likelihood $p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}_k)$ as the weights are no longer shared. But this rises more problem: (1) it requires more samples, and it might be impossible in real-world scenarios. (2) it requires samples generated by each memory pattern $\boldsymbol{\xi}_k$, as each individual weight \mathbf{w}_k is optimized by samples involving $\boldsymbol{\xi}_k$.

Additionally, the adaptive similarities are a family of similarity measures that can fit to the variant distribution through sampling, it is not solely Equation 4. When proving the optimal correct retrieval for noisy, masked, and biased variants, we propose more adaptive similarity as theoretical tools. We wrote Equation 4 in the main text simply because it is the most effective ones for retrieving under noisy, masked, biased, and mixed settings, and it requires minimum trainable weight, as examined in ablation study (Appendix A.4.6).

One more thing is that the priori $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$ are often ignored in this work. Well, it should not be ignored in all cases. However, we can use a bias term $\mathbf{b} \in \mathbb{R}^N$ and use \mathbf{b}_k capture the occurrance of pattern $\boldsymbol{\xi}_k$. By assuming similarity is a logits of the posteriori $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$ (e.g., $\log p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$), then we can see that $\arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \{\log p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})\} = \arg \max_{\boldsymbol{\xi}' \in \boldsymbol{\Xi}} \{\log p_{\mathcal{V}}(\mathbf{x}|\boldsymbol{\xi}') + \log p_{\mathcal{V}}(\boldsymbol{\xi}')\}$, and we set pass $\mathbf{s}(\boldsymbol{\Xi}, \mathbf{x}) + \mathbf{b}$ to the separation function so that \mathbf{b} handles the term $\log p_{\mathcal{V}}(\boldsymbol{\xi})$.

A.4 EXPERIMENTS

A.4.1 BASELINES AND METRICS

For memory retrieval test, we compare our Adaptive Hopfield network A-Hop against: Modern Hopfield network M-Hop (Ramsauer et al., 2021); Universal Hopfield network U-Hop (Millidge et al., 2022) with $sim(\boldsymbol{\xi}, \mathbf{x}) = -\|\boldsymbol{\xi} - \mathbf{x}\|_1$ and $sep(\cdot) = arg max(\cdot)$ as they report a leading

performance for such configuration when masking out half of the dimensions in memory patterns (similar to masked variant, Definition 5); Kernelized Hopfield network K-Hop (Wu et al., 2024a) with the kernel optimized by separation loss proposed by Wu et al. (2024a); Kernelized Hopfield network K₂-Hop with the kernel optimized by loss defined by variant distribution (Equation 5) rather than the separation loss; and Multi-Layer Perceptrons MLP that has 4 layers and a input dimensionality d, output dimensionality N, trying to fit $p_{\mathcal{V}}(\boldsymbol{\xi}|\mathbf{x})$ but unsatisfactory. We estimate the empirical retrieval accuracy of each model (Definition 9), and the mean squared error $\mathbb{E}_{(\boldsymbol{\xi},\mathbf{x})\sim\mathcal{V}(\boldsymbol{\Xi})}\left[(\mathcal{T}(\mathbf{x})-\boldsymbol{\xi})^{\top}(\mathcal{T}(\mathbf{x})-\boldsymbol{\xi})\right]$. We wrote a generator that can generate noisy, masked, biased, and mixed variants.

For tabular classification test, we compare the A-Hop with a memory-based classifier (Appendix A.4.4) with: M-Hop uses the same classifier framework but the unlearnable similarity function; K_2 -Hop that uses the same classifier framework but a different similarity function, and its kernel function is optimized by the classification loss; Extremely Randomized Trees (Geurts et al., 2006), or Extra Trees, a tree-based classifier; Random Forest (Breiman, 2001), yes, the famous Random Forest classifier; AdaBoost (Freund & Schapire, 1997), a classic boosting classifier; and XGBoost (Chen & Guestrin, 2016), an enhanced Gradient Boosting Decision Tree. All models are tuned with a 5-fold cross-validation on the training set. We measure the test accuracy for the dataset with a positive sample rate 0.2 < % pos < 0.8, and the ROC-AUC score otherwise, following the settings in Wang et al. (2025).

For image classification task, we follow the settings in Wu et al. (2024a), where they replaced the attention component with a HopfieldLayer (Ramsauer et al., 2021). We experiment by integrating A-Hop, M-Hop, and K-Hop into the HopfieldLayer. Similarly, we follow the settings in Ramsauer et al. (2021); Hu et al. (2023), where they use HopfieldPooling for multiple instance learning. We integrate A-Hop, M-Hop, and K-Hop to HopfieldPooling, and run a 5-fold cross-validation to report the mean ROC-AUC of all folds as the result.

For all experiments, we report the results with the mean and standard deviation of five runs.

A.4.2 DATASETS

We used a total of 12 datasets to assess the performance of A-Hop on tasks including tabular classification (Adult, Bank, Vaccine, Purchase, and Heart), and image classification (CIFAR 10, CIFAR 100, and Tiny ImageNet), and multiple instance learning (Tiger, Fox, Elephant, UCSB).

- **Adult** (Becker & Kohavi, 1996) The prediction task for this dataset is to classify individuals' income levels as either above or below \$50,000 annually. The data was extracted by Barry Becker from the 1994 Census database.
- **Bank** (Moro et al., 2014) To predict the success of a term deposit subscription, this dataset records the outcomes of telemarketing campaigns from a banking institution in Portugal.
- **Vaccine** (Bull et al., 2016; DrivenData, 2019) We use this dataset from a DrivenData competition to predict if a person received a seasonal flu vaccine. The data consists of 26,707 survey responses detailing 36 behavioral and personal attributes.
- **Purchase** (Sakar & Kastro, 2018) The objective with this dataset is to forecast the online shopping intentions of visitors to an e-commerce website, determining whether a user will proceed with a purchase.
- **Heart** (Kaggle, 2018) This dataset facilitates the prediction of cardiovascular disease presence. It contains health-related data from 70,000 patients, as provided by Ulianova.
- **CIFAR 10** (Krizhevsky et al., 2009) is a classic image recognition dataset consisting of 60,000 32×32 color images in 10 classes, with 6,000 images per class.
- **CIFAR 100** (Krizhevsky et al., 2009) is a more challenging version of CIFAR 10, containing the same number of images but split into 100 fine-grained classes.
- **Tiny ImageNet** (Le & Yang, 2015) is a subset of the ImageNet dataset designed for educational purposes. It contains 100,000 images from 200 classes, downsized to 64×64 pixels.
- **Tiger, Fox, Elephant** (Deng et al., 2009) These are specific class subsets extracted from the large-scale ImageNet database, often used for fine-grained image classification tasks.

UCSB (Gelasca et al., 2008) This dataset is composed of 58 breast cancer histopathology images stained with Hematoxylin and Eosin (H&E). The primary challenge it presents is the accurate segmentation of individual cells from the complex tissue background, which is a critical precursor to classifying cells as benign or malignant.

A.4.3 MEMORY RETRIEVAL

We first introduce the intensity of variant setting. A triplet $(d_{\text{mask}}, d_{\text{noise}}, d_{\text{bias}}) \in [0, 1]^3$ is used to describe a variant setting. This mean that we will first add a Gaussian noise $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, d_{\text{noise}}\mathbf{I})$ to a certain memory pattern $\boldsymbol{\xi} \in \boldsymbol{\Xi}$ obtaining $\mathbf{x} \leftarrow \boldsymbol{\xi} + \mathbf{n}$. Then, we will choose $d \cdot d_{\text{mask}}$ indices, and set these indices of \mathbf{x} to random numbers choosen uniformly from [-1,1]. Finally, we will add a bias \mathbf{d} to $\mathbf{x} \leftarrow \mathbf{x} + \mathbf{d}$ with $\mathbf{d}_i = \mathbf{s}_i \cdot d_{\text{bias}}$, where \mathbf{s}_i is a random sign sampled uniformly from $\{-1,+1\}$. Then, the generator return $(\boldsymbol{\xi}, \mathbf{x})$ as a sample of $\mathcal{V}(\boldsymbol{\Xi})$. Therefore, different triplet describle different variant settings, and thus, corresponds to different variant distribution $\mathcal{V}(\boldsymbol{\Xi})$.

In the memory retrieval task, we use the retrieval dynamics written in Equation 4. To learn the weight w's and β 's we use optimizer Adam for 200 epoches, and use a learning rate 0.1 in all settings. However, we argue that number of epoches (N_epoch) and learning rate (1r) should be tuned when applying A-Hop to other memory retrieval settings.

For K-Hop and K₂-Hop, we tuned them carefully. For the original K-Hop, it optimize a separation loss, and we try to make it as small as possible. However, the retrieval accuracy of K-Hop is not satisfactory, and it is reasonable since it is not optimized for correct retrieval, but for ϵ -retrieval. We found that the domain of the memory pattern in their experiments is $[0,1]^d$, and it is harmful to dot product-based methods (think of using only $1/2^d$ of the space) as dot product highly relies on signs. Therefore, for fair comparision, we sample the random vectors uniformly from $[-1,1]^d$ and change the domain of pixels in MNIST images to [-1,1].

A.4.4 TABULAR CLASSIFICATION

We develop a memory-based model for tabular classification that takes advantage of the excellent memory retrieval effectiveness of associative memories. The insight is that classification is hierarchical, and we divide instances into *cases*, and further classify cases into the final class. For instance, there are type I diabetes and type II diabetes, where each type here resembles the idea of cases. Another samples is that cats has plenty of breeds, and associative memory can capture specific idea of orange cat, or blue cat, while they have a hard time figuring out the generalize idea of cats. So, we let associative memory match instances into cases by choosing some instances in the training set as the representatives of cases, or use K-means cluster to produce such representatives, and pass the retrieval probability to a multi-layer perceptron (MLP) for classifying cases. That is, the model can be represented as $\mathbf{y} = \text{MLP}(\text{LayerNorm}(\sin(\mathbf{\Xi}, \mathbf{x})))$, and this can be trained on a conventional machine learning fashion by minimizing classification loss:

$$\mathcal{L}(\mathcal{D}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[(\mathbf{y} - \hat{\mathbf{y}})^2]$$
(9)

We no longer tune weights in adaptive memories using loss defined in Equation 5, as they can be tuned using the classification loss when participanting into going forward in the network.

We tuned the hyperparameters by grid searching on the training data:

Table 6: Hyperparameters tuned in tabular classification

Name	Domain
N_epoch batch_size	
init	Cluster, No cluster

A.4.5 IMAGE CLASSIFICATION AND MULTIPLE INSTANCE LEARNING

For image classification, we follow the settings in Wu et al. (2024a), and place the adaptive similarity to the Hopfield layer inside a image Transformer. For M-Hop and K-Hop, we use the hyperparame-

ter suggested in Wu et al. (2024a), and find the optimal number of epoch and learning rate for A-Hop via grid search. We run five iterations of separate loss optimization for A-Hop before testing.

Table 7: Hyperparameters tuned in image classification

Name	Domain
N_epoch	$\{25, 40, 50\}$
lr	$\{10^{-3}, 10^{-4}\}$

For multiple instance learning, we follow the settings in Hu et al. (2023), and use HopfieldPooling as the backbone. Similarly, we place the adaptive similarity in the core Hopfield component replacing the fixed dot product. All experiments are run on a 5-fold validation, and the ROC-AUC scored is taken as the mean of all folds.

Table 8: Hyperparameters tuned in multiple instance learning

Name	Domain
N_epoch lr lr_decay	

A.4.6 ABLATION STUDY

We conduct four different ablation studies to see the effective of different components.

In the first experiment (Table 9), we tested if sorting the dimension-wise similarity vector \mathbf{q} and the upper-right triangle matrix \mathbf{U} is needed. The results shows that both of them are necessary for high retrieval accuracy.

Table 9: Retrieval accuracy (\uparrow) and error (\downarrow) between unsorted and sorted q. Each cell contains the mean accuracy or error with standard deviation in a smaller font. Results of the best-performing model are **bolded**.

Conditions		Synthetic	(d = 0.4)	Synthetic ($d = 0.5$)		
q sorted?	use \mathbf{U} ?	Accuracy	Error	Accuracy	Error	
X	X	.5172±.034	.1900±.017	.2094±.022	.2658±.005	
×	\checkmark	$.5444 \pm .007$	$.1665 {\pm} .007$.1888±.015	$.2738 \pm .003$	
✓	X	$.6928 \pm .034$	$.1173 \pm .010$.3374±.025	$.2324 \pm .006$	
\checkmark	\checkmark	.7280 ±.034	.1033 ±.011	.3634 ±.040	.2207 \pm .011	

Next, we look for a better matrix \mathbf{U} (Table 10), which is the core of similarity measure. Our results show that the upper-right triangle structure of \mathbf{U} is optimal, while making \mathbf{U} learnable and initialize it with the upper-right triangle yields the best result, but we does not adopt learnable \mathbf{U} in other experiments to keep minimum learnable parameters. Another finding is that a randomized matrix \mathbf{U} is better than not having \mathbf{U} (when $\mathbf{U} = \mathbf{I}$). Therefore, along with Table 9, we find that the footprint structure is essential to adaptive similarities.

Table 10: Retrieval accuracy (\uparrow) and error (\downarrow) between different configuration of U. Each cell contains the mean accuracy or error with standard deviation in a smaller font. Results of the best-performing model are **bolded**, and the second are <u>underlined</u>.

Conditions		Synthetic	(d = 0.4)	Synthetic ($d = 0.5$)	
U learnable?	initialization?	Accuracy	Error	Accuracy	Error
×	Random	.6928±.015	$.1147 \pm .005$.3340±.019	$.2305 \pm .005$
×	I	$.6802 \pm .013$	$.1194 \pm .004$.3458±.025	$.2298 \pm .008$
X	${f U}$	$.7242 \pm .016$	$.1056 \pm .007$	<u>.3600</u> ±.024	$.2272 \pm .005$
✓	Random	$.7176 {\scriptstyle \pm .027}$	$.1087 {\pm} .007$.3414±.023	$.2292 {\pm} .006$
✓	I	$.7114 \pm .019$	$.1096 \pm .006$.3528±.021	$.2270 \pm .005$
\checkmark	\mathbf{U}	.7280 ±.034	.1033 ±.011	.3634 ±.040	.2207 ±.011

We also tested the effect of different footprint with different base similarity (see Table 11). Results shows that two footprints (dis and dot) is always better than one alone, and we found that $ftpt_{dis}$ performs better than $ftpt_{dot}$ in this setting. However, even though retrieval accuracy degrades by removing one of the footprint, using $ftpt_{dis}$ or $ftpt_{dot}$ only is still better than other Hopfield networks.

Table 11: Retrieval accuracy (\uparrow) and error (\downarrow) between different usage of footprint. Each cell contains the mean accuracy or error with standard deviation in a smaller font. Results of the best-performing model are **bolded**.

Conditions		Synthetic	(d = 0.4)	Synthetic ($d = 0.5$)		
use $\mathrm{ftpt}_{\mathrm{dis}}$?	use $ftpt_{dot}$?	Accuracy	Error	Accuracy	Error	
×	✓	$.5926 \pm .016$	$.1517 \pm .005$.2520±.027	$.2515 \pm .004$	
\checkmark	X	$.6458 \pm .042$	$.1317 \pm .011$.3286±.023	$.2326 {\pm}.007$	
\checkmark	✓	.7242 ±.016	$\boldsymbol{.1056} {\pm .007}$.3600 ±.024	.2272 ±.005	

Finally, we tested the number of samples needed for adaptive similarity to mimic the variant distribution. It looks like training on only 512 samples provides a good enough adaptive similarity for 2048 64-dimension memory patterns.

Table 12: Retrieval accuracy (\uparrow) and error (\downarrow) between different number of samples provided for learning. Each cell contains the mean accuracy or error with standard deviation in a smaller font. Results of the best-performing model are **bolded**.

Conditions	Synthetic ($d = 0.4$)		Synthetic ($d = 0.5$)	
number of samples?	Accuracy	Error	Accuracy	Error
512	.7204±.029	.1077±.010	.3541±.028	.2311±.004
∞	.7242 ±.016	.1056 ±.007	.3600 ±.024	.2272 ±.005