

# Modern Methods in Associative Memory

Dmitry Krotov<sup>1,2</sup>, Benjamin Hoover<sup>1,3</sup>, Parikshit Ram<sup>1</sup>, Bao Pham<sup>1,4</sup>

<sup>1</sup>IBM Research, <sup>2</sup>MIT, <sup>3</sup>Georgia Tech, <sup>4</sup>RPI

**Date:** July 14, 2025

**Tutorial:** ICML 2025, *Vancouver, BC, Canada*

**Website:** <https://tutorial.amemory.net>

## Abstract

Associative Memories like the famous Hopfield Networks are elegant models for describing fully recurrent neural networks whose fundamental job is to store and retrieve information. In the past few years they experienced a surge of interest due to novel theoretical results pertaining to their information storage capabilities, and their relationship with SOTA AI architectures, such as Transformers and Diffusion Models. These connections open up possibilities for interpreting the computation of traditional AI networks through the theoretical lens of Associative Memories. Additionally, novel Lagrangian formulations of these networks make it possible to design powerful distributed models that learn useful representations and inform the design of novel architectures. This tutorial provides an approachable introduction to Associative Memories, emphasizing the modern language and methods used in this area of research, with practical hands-on mathematical derivations and coding notebooks.

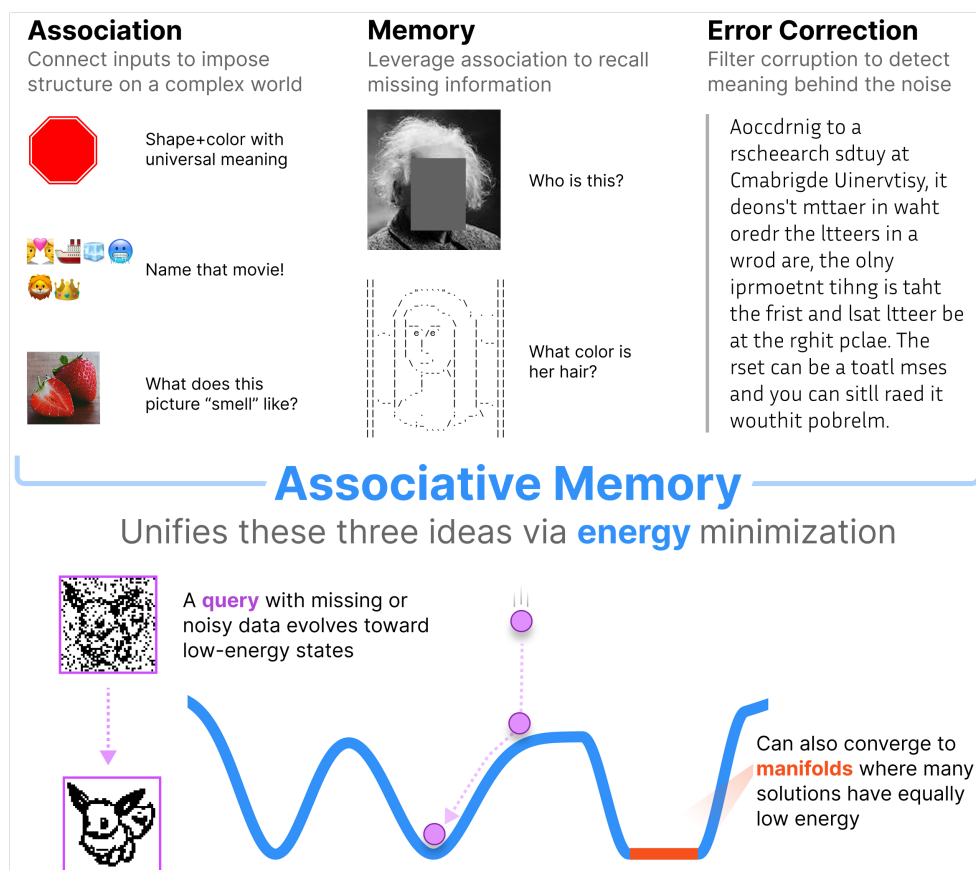
# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Dense Associative Memory: Discrete State Vector</b>	<b>5</b>
2.1	Information Storage Capacity . . . . .	6
2.2	Limiting Cases . . . . .	9
2.3	General Dense Associative Memory with Binary State Variables . . . . .	10
<b>3</b>	<b>General Dense Associative Memory</b>	<b>12</b>
3.1	Building Blocks of AMs with Modular Energies . . . . .	12
3.2	Dynamical Neurons and their Lagrangians . . . . .	14
3.3	Hypersynapses . . . . .	15
3.4	Energy Descent Dynamics . . . . .	17
3.5	Implementing AMs . . . . .	17
3.5.1	Energy Transformer Block . . . . .	18
3.6	Bridging Energy Minimization and Feedforward Prediction . . . . .	23
<b>4</b>	<b>Failure of Memory and Generative AI</b>	<b>25</b>
4.1	Diffusion Models . . . . .	26
4.2	Diffusion Models from Associative Memories . . . . .	27
4.3	Memorization - Spurious - Generalization Transition . . . . .	30
<b>5</b>	<b>Associative Memory: A Machine Learning Model</b>	<b>33</b>
5.1	Machine Learning Modeling . . . . .	33
5.2	Associative Memory Network as a Model . . . . .	36
5.2.1	Memory Capacity and Expressivity . . . . .	38
5.2.2	Supervised Learning . . . . .	39
5.2.3	Nonparametric vs Parametric Models . . . . .	40
5.3	Clustering . . . . .	40
5.3.1	Euclidean Clustering . . . . .	41
5.3.2	Deep Clustering . . . . .	44
5.4	Kernel Machines . . . . .	46
5.4.1	Random Features . . . . .	47
5.4.2	Novel Energy Functions . . . . .	49
<b>6</b>	<b>Conclusion</b>	<b>54</b>

# Chapter 1

## Introduction

Associative Memory (AM) is a core concept in psychology responsible for linking related items [1]. For instance, if one is shown an image of a strawberry, it is likely that they can recall the smell and taste of this fruit; or, in the case of an image of a person, an acquaintance of them would be able to name them, see Fig. (1.1) for the demonstration of AM. These are examples of input-output pairs that are associated in our memory, where prompting for an element of the pair results in content-addressable retrieval of the other element.



**Figure 1.1:** The form of Associative Memory discussed in this tutorial uses an energy function to unify three important aspects of human cognition: association, memory, and error correction. We are capable of associating images, sights, sounds, smells, and symbols with each other. These associations allow us to retrieve memories using partial or corrupted information, making it a content-addressable memory with error-correction capabilities. The functionality of Associative Memory is modeled by an energy function, where low values of the energy correspond to stored memories and constitute the most likely states of the system.

Another important aspect of Associative Memory is the notion of error correction. You can easily read the text in Fig. (1.1) without much difficulty, despite the fact that almost no words in that paragraph are proper English words. The reason why you are able to comprehend this text is because there are powerful error correction mechanisms that are constantly working in your brain that associate imperfect inputs with the proper semantic meaning of individual words. The same applies to the example above: the image of the strawberry can be presented with all kinds of distortions and imperfections. Despite all that variability in the input, Associative Memory manages to link those inputs with the proper smell and taste.

Thus, **Associative Memory is a content addressable information storage system that is capable of error correction.**

Associative Memory plays a major role in the history of AI. Following the 1943 model of artificial neuron by McCulloch and Pitts [2], and a body of work [3] on artificial neural networks (ANNs) – Perceptrons – by Frank Rosenblatt in the 1950-1960s, the public community at large has been extremely enthusiastic about the future. Popular media outlets from that period promised that the Perceptron “will be able to walk, talk, see, write, reproduce itself and be conscious of its existence” [4], similar to what we read in popular press about AI today. However, in 1969, Minsky and Papert demonstrated that simple Perceptrons (without hidden layers) could not compute even the simplest logical gates, e.g., XOR [5]. The public perception of this result led to the drop of enthusiasm in ANNs. Most of the computer science community at that time left the field of ANNs — triggering what historians of science later called the “AI Winter” [6].

John Hopfield’s seminal paper of 1982 [7] on what is now called the Hopfield network of Associative Memory was the major driving force that ended that period. In his paper, Hopfield connected computational aspects of Associative Memory with collective properties of Ising [8] magnets in condensed matter physics, which were a “hot topic” at the time. Specifically, Hopfield posed a simple quantifiable problem: Given a network of  $D$  neurons, how much information (or memories) can such a network store and retrieve? Content-addressable Associative Memory retrieval was a sufficiently non-trivial problem to illustrate the potential of ANNs’ computational abilities. At the same time, it was simple enough to be analytically solvable using powerful methods developed in statistical physics. The confluence of these two aspects created a “harmonic oscillator” level abstraction for ANN computation, and set the grounds for many extensions and generalizations that followed.

Associative memory has been a prominent theme of ANN research in the 1960s-1980s. A highly incomplete and subjective list of seminal papers from that period includes: Anderson [9], Willshaw, Buneman, Longuet-Higgins [10], Amari [11; 12], Cohen and Grossberg [13], Hopfield [14], Amit, Gutfreund, Sompolinsky [15], and many others.

The main focus of this tutorial is on the **Energy-based Associative Memories**. This is the class of ANNs which are recurrent neural networks that can be described by a state vector, which evolves in time according to some non-linear rule. This state vector can be either continuous or discrete. The update rule can be written either in terms of continuous time (differential equation), or discrete set of update steps (which we will usually treat as a discretization of that differential equation). Finally, there exist multiple options for how one updates the state vector.

The most common choices are: *synchronous* — all compute elements of the state vector are updated simultaneously, or *asynchronous* — at any given time a subset of all elements of the state vector are updated, i.e., a random element of the state vector is updated while the remaining elements are kept intact. For the purposes of this tutorial, we will most work with continuous states, continuous time, and synchronous updates. Thus, the state vector of the network  $\mathbf{x} \in \mathbb{R}^D$ , which has individual elements  $x_i$  (index  $i$  runs from 1 to  $D$ ), evolves according to the following differential equation:

$$\frac{dx_i}{dt} = f_i(\mathbf{x}, t) \quad (1.1)$$

where the functions  $f_i(\mathbf{x}, t)$  represent the vector field that defines the dynamics. We will refer to individual elements of this vector as “neurons” although in many situations these variables may describe a different biological structure, e.g., an astrocyte or their processes (long tentacles originating from the astrocyte’s cell body).

A general system of coupled non-linear differential equations may have many complex behaviors: fixed points, limit cycles, strange attractors, or chaotic behavior. Energy-based AMs are a special subclass of general systems (1.1) that have the notion of an energy function (sometimes also referred to as a Lyapunov function). You can think of the temporal evolution of the state vector as a ball rolling downhill in a sophisticated energy landscape as seen in Fig. (1.1). The energy is bounded from below, and the ball is only allowed to move in a way that decreases its energy. Because of these restrictions, eventually, the ball must either stop at one of the local minima or reach a manifold that corresponds to flat energy. In the latter case, the ball may continue to move along that manifold as long as the energy does not increase.

The local minima of the energy (which can be point-like attractors — zero-dimensional manifolds — or alternatively, manifolds of higher dimension) are called **memories**. The process of shaping the energy landscape corresponds to writing information into the AM network or learning. The dynamical trajectory of energy descent, illustrated by Eq. (1.1), corresponds to memory recall, or inference. Association happens between the initial state of the network  $\mathbf{x}(t = 0)$  and the final asymptotic state of the network  $\mathbf{x}(t \rightarrow \infty)$ . Finally, the asymptotic states of the dynamics are typically stable (unless they lie on the flat portions of the energy landscape). Intuitively, this means that small perturbations that do not push the state vector outside the basin of the fixed point’s attraction gets auto-corrected by the network itself. For these reasons, this network is an AM system.

In some settings, memories in this system may correspond to individual instances of the training data. Alternatively, they may correspond to emergent attracting manifolds that are shaped by the learning algorithm (e.g., backpropagation [16; 17], Hebbian learning [18], contrastive training [19; 20], etc.). In the latter case, the memories do not typically correspond to individual instances of the training data, but rather describe consolidated memories — “knowledge” — that the network acquires through synergy of AM architecture, a specific learning algorithm, and the choice of training data.

Intuitively, you can think about the initial state  $\mathbf{x}(t = 0)$  as a “question” that you ask the neural network. This question positions the state vector at some high-energy location on the energy landscape. The network will perform computation by moving that state to a local minimum

(or a metastable state) — the process of “thinking.” Once the local minimum is reached, the computation stops and the network’s state does not evolve in time anymore. You can read out that final state  $\mathbf{x}(t \rightarrow \infty)$  and convert it to the answer to the posed question. Importantly, this computation is very different from conventional feed-forward architectures, e.g., feed-forward convolutional neural networks, transformers, or large language models without chain of thought. These conventional architectures are described by a computational graph with a finite number of steps. This means that if the network has 10 layers, it must produce some kind of an answer to the question after exactly 10 steps. This happens regardless of the complexity of the posed question. AM architectures are very different. They can dynamically adapt the computational graph based on the complexity of the posed question. For simple questions the network may produce an answer in 5 steps, but for more complicated questions the network may need to “think” longer.

Finally, because of the network’s energy-based architecture, the final answer is **asymptotically stable**. This means that once the computation or “thinking” stopped and the network converged to an answer, the precise timing of the output’s readout doesn’t matter. Assuming that the readout time  $T$  is large enough, we can use  $\mathbf{x}(t = T)$  or  $\mathbf{x}(t = T + 0.5 \text{ seconds})$  as the final answer, and the two must be identical. This property of asymptotic stability makes AM framework extremely appealing for neuromorphic devices, where hardware imperfections may prevent the ability to read the network’s state at a precise timing.

In the past few years there has been significant advances in the field of AMs. These advances pertain to the development of **Dense Associative Memories** or DenseAMs [21]. They are flexible energy-based AM architectures that are capable of storing large amounts of information, enable incorporation of many useful inductive biases (e.g., convolutions [22], attention [23], etc.) in their architecture, and have mathematically controllable properties of emergent local minima. DenseAM ideas have triggered a large amount of innovative ideas about the potential use cases of AMs and we believe they will enable a new frontier for AM research [24]. In this tutorial, we will cover many of these new developments from both the theoretical perspective and practical implementations. The tutorial is supplemented with a collection of notebooks and suggested problems that the readers can explore on their own to better understand the core ideas and methods, and to get hands-on experience coding an AM network suitable for their own use case. We intentionally designed these problems so that they are simple but still illustrate a useful mathematical concept and the wonderful idea of AM. We hope that you enjoy this learning experience.

## Chapter 2

# Dense Associative Memory: Discrete State Vector

This chapter introduces a popular class of AMs — Dense Associative Memory (DenseAM). This family of models is a generalization of celebrated Hopfield networks. While Hopfield networks are very elegant mathematical models that satisfy all of the AM requirements, they are known to have a very small information storage capacity. As a result, DenseAMs are specifically designed to retain all of the benefits within Hopfield networks, but rectify their small information storage issue [21; 25].

As discussed earlier, AMs can be formulated both in discrete and continuous variables, and in discrete or continuous time. In this chapter, we focus on DenseAMs with discrete state vector, and discrete asynchronous updates. Specifically, we will be working with a set of discrete variables  $\sigma_i = \{\pm 1\}$ , index  $i = 1, \dots, D$  which compose a state vector  $\boldsymbol{\sigma}$ . In addition to that, the network will have  $K$  memory vectors  $\boldsymbol{\xi}^\mu$  with index  $\mu = 1, \dots, K$ . Each memory is a  $D$ -dimensional vector with individual elements denoted by  $\xi_i^\mu$ .

The energy function is defined as:

$$E = - \sum_{\mu=1}^K F\left(\sum_{i=1}^D \xi_i^\mu \sigma_i\right). \quad (2.1)$$

The goal of the network is to start at some initial state  $\sigma_i^{(t=0)}$ , which typically corresponds to a high-energy state, and lower the energy by flipping the elements of the state vector. The dynamics of flipping stops when no further single element flip can reduce the energy. At that point, the network has reached a local minimum of the energy. As usual, we will refer to the individual elements of the state vector as neurons or spins.

In order to formalize this intuitive dynamical equation, pick a single neuron and define its state

at the next iteration as:

$$\begin{aligned}
\sigma_i^{(t+1)} &= \text{Sign} \left[ E(\sigma_i = -1, \sigma_{j \neq i} = \sigma_j^{(t)}) - E(\sigma_i = +1, \sigma_{j \neq i} = \sigma_j^{(t)}) \right] \\
&= \text{Sign} \left[ \sum_{\mu=1}^K F\left(\xi_i^\mu + \sum_{j \neq i}^D \xi_j^\mu \sigma_j^{(t)}\right) - \sum_{\mu=1}^K F\left(-\xi_i^\mu + \sum_{j \neq i}^D \xi_j^\mu \sigma_j^{(t)}\right) \right] \\
&\approx \text{Sign} \left[ \sum_{\mu=1}^K 2 \xi_i^\mu F'\left(\sum_{j \neq i}^D \xi_j^\mu \sigma_j^{(t)}\right) + \text{higher order subleading terms} \right].
\end{aligned} \tag{2.2}$$

This update rule compares the energies of two states:  $\sigma_i = -1$  with states of all other neurons clamped to their current values, and  $\sigma_i = +1$  with all the other neurons clamped. The  $\text{Sign}[\cdot]$  function assigns the state of the  $i^{\text{th}}$  neuron to the one corresponding to the lowest energy among these two possibilities. Finally, in the last line of Eq. (2.2) we have used the Taylor series to expand the function  $F(x + \varepsilon) \approx F(x) + \varepsilon F'(x) + \text{higher order terms}$ . It is legitimate to terminate the expansion after the first term since  $\varepsilon = \pm 1$  is much smaller than the overlap between the clamped part of the state vector and the memories.

This update rule is typically written in the following form:

$$\sigma_i^{(t+1)} = \text{Sign} \left[ \sum_{\mu=1}^K \xi_i^\mu f\left(\sum_{j \neq i}^D \xi_j^\mu \sigma_j^{(t)}\right) \right], \tag{2.3}$$

where we dropped the factor of 2 in the argument of the sign function (it doesn't play any role there) and introduced an activation function  $f(\cdot) = F'(\cdot)$ , which is a derivative of the function  $F(\cdot)$  defining the energy.

The energy function (2.1) is a finite sum of smooth functions (we assume that the function  $F(\cdot)$  does not have singularities - infinite values for finite arguments) that depend on the finite number of discrete variables. Thus, the energy is finite and bounded from below. Additionally, the dynamical equations (2.2) and (2.3) decrease the value of energy at each iteration. Thus, if we keep applying these update equations to the state vector for a long time, eventually the system will reach a steady state – no single neuron flip can further reduce the energy.

## 2.1 Information Storage Capacity

How many memories or local minima can such a system store and successfully retrieve? The network, specified by Eq. (2.1), can be defined for any number  $K$  of memories. But, it turns out, if you pack too many of such vectors inside the  $D$ -dimensional discrete space, the local minima of the energy will no longer correspond to the stored patterns. In what follows we will compute the largest value of  $K$  that permits successful remembering of the stored patterns.

In general, this maximal value  $K^{\max}$  will depend on the specific choices for the stored memories. We will derive a statistical scaling law for this memory capacity assuming that the patterns are

drawn at random from the following distribution:

$$\xi_i^\mu = \begin{cases} +1, & \text{with probability } \frac{1}{2} \\ -1, & \text{with probability } \frac{1}{2} \end{cases} \quad (2.4)$$

With this distribution, it is easy to compute the correlation functions for these variables. The one-point and two-point correlation functions are given by:

$$\langle \xi_i^\mu \rangle = 0, \quad \langle \xi_i^\mu \xi_j^\nu \rangle = \delta^{\mu\nu} \delta_{ij} \quad (2.5)$$

In order to quantify the information storage capacity of this network we will use the following trick. We will initialize the network in the state corresponding to one of the memories, say  $\xi_i^1$ , and let it evolve in time according to the update rule. If the pattern  $\xi_i^1$  corresponds to a local minimum, that state must be stable. In other words, the dynamics should not change that initial state. Mathematically, this means that

$$\begin{aligned} \sigma_i^{(t+1)} &= \text{Sign} \left[ \xi_i^1 f \left( \sum_{j \neq i}^D \xi_j^1 \xi_j^1 \right) + \sum_{\mu=2}^K \xi_i^\mu f \left( \sum_{j \neq i}^D \xi_j^\mu \xi_j^1 \right) \right] \\ &= \text{Sign} \left[ \underbrace{\xi_i^1 f(D-1)}_{\text{signal}} + \underbrace{\sum_{\mu=2}^K \xi_i^\mu f \left( \sum_{j \neq i}^D \xi_j^\mu \xi_j^1 \right)}_{\text{noise}} \right] \stackrel{?}{=} \xi_i^1 \end{aligned} \quad (2.6)$$

### Derivation of the generating function

It is helpful to introduce a new variable

$$\Xi = \sum_{j=2}^D \xi_j^1 \quad (2.7)$$

and compute the generating function defined as a statistical average of the exponent of that variable

$$M(\tau) = \langle e^{\tau \Xi} \rangle \quad (2.8)$$

Since  $\xi_j^1$  are independent for different indices  $j$ , the statistical average can be factorized and computed explicitly

$$M(\tau) = \frac{1}{2^{D-1}} \sum_{\xi_2=\pm 1} \sum_{\xi_3=\pm 1} \dots \sum_{\xi_D=\pm 1} e^{\tau \xi_2} e^{\tau \xi_3} \dots e^{\tau \xi_D} = \cosh(\tau)^{D-1} \quad (2.9)$$

All correlation function can be computed by taking derivatives of the generating function. For instance

$$\langle \Xi^{2p} \rangle = \left. \frac{\partial^{2p} M}{\partial \tau^{2p}} \right|_{\tau=0} = (2p-1)!! D^p \quad (2.10)$$

Assuming that the function  $f(\cdot)$  is non-negative, the signal term pushes the argument of the Sign function towards aligning it with the desired pattern  $\xi_i^1$ . The noise term generally pushes that argument away from the desired pattern and in some situations may outweigh the signal term. Below, we will compute the characteristic magnitude of the noise term and determine when it becomes dominant and destroys the stability of the target memory. Specifically, we can compute the mean and variance of the noise term. The mean

$$\langle \text{noise} \rangle = \left\langle \sum_{\mu=2}^K \xi_i^\mu f\left(\sum_{j \neq i}^D \xi_j^\mu \xi_j^1\right) \right\rangle = 0 \quad (2.11)$$

is equal to zero since index  $i$  appears only once in the correlator, see Eq. (2.5). The variance is given by

$$\begin{aligned} \langle \text{noise}^2 \rangle &= \left\langle \sum_{\mu=2}^K \xi_i^\mu f\left(\sum_{j \neq i}^D \xi_j^\mu \xi_j^1\right) \sum_{\lambda=2}^K \xi_i^\lambda f\left(\sum_{k \neq i}^D \xi_k^\lambda \xi_k^1\right) \right\rangle \\ &= \sum_{\mu=2}^K \left\langle f\left(\sum_{j \neq i}^D \xi_j^\mu \xi_j^1\right) f\left(\sum_{k \neq i}^D \xi_k^\mu \xi_k^1\right) \right\rangle \stackrel{\text{i.d.}}{=} \sum_{\mu=2}^K \left\langle f\left(\sum_{j \neq i}^D \xi_j^\mu\right) f\left(\sum_{k \neq i}^D \xi_k^\mu\right) \right\rangle \quad (2.12) \\ &= (K-1) \left\langle f\left(\sum_{j \neq i}^D \xi_j^\mu\right)^2 \right\rangle, \end{aligned}$$

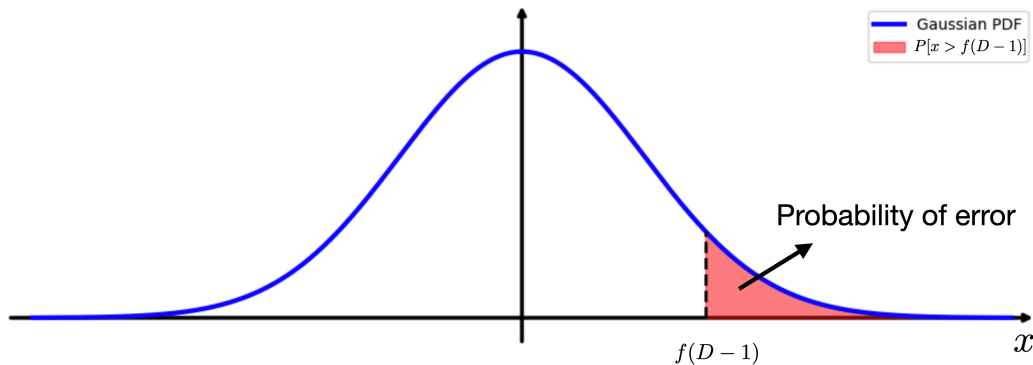
where we used that  $\langle \xi_i^\mu \xi_i^\lambda \rangle = \delta^{\mu\lambda}$  and the property that in distribution  $\xi_j^\mu \xi_j^1 \stackrel{\text{i.d.}}{=} \xi_j^\mu$ .

Now, it is instructive to restrict our calculation to the class of power energy functions so that

$$F(\cdot) = \frac{1}{n}(\cdot)^n, \quad f(\cdot) = (\cdot)^{n-1}, \quad \text{where } n \text{ is an integer.} \quad (2.13)$$

In this case, the variance of the noise can be computed exactly (through the generating function) and is equal to<sup>1</sup> [26]:

$$\Sigma^2 = \langle \text{noise}^2 \rangle = (2n-3)!! K D^{n-1}. \quad (2.14)$$



**Figure 2.1:** Gaussian probability distribution function. Shaded area indicates the probability of an error or spin flip.

Now we are ready to compute the probability of an error. The noise term in Eq. (2.6) is a sum

<sup>1</sup>We assume that  $K$  is large so that  $K-1 \approx K$ .

over many independent random variables. When  $K$  and  $D$  are large, this noise term behaves approximately as a Gaussian random variable. When the sign of the noise term is the same as the sign of the signal, the noise term pushes the update in the right direction and does not cause issues. The problem arises when the noise is large and its sign is opposite to that of the signal. In this situation, it is possible that the noise can outweigh the signal and flip the spin of interest. The probability of this event is given by the area under a Gaussian distribution, as shown in Fig. (2.1):

$$P(\text{error}) = \int_{f(D-1)}^{\infty} \frac{dx}{\sqrt{2\pi}\Sigma^2} e^{-\frac{x^2}{2\Sigma^2}} = \int_{\frac{f(D-1)}{\Sigma}}^{\infty} \frac{dy}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} = g\left(\frac{f(D-1)}{\Sigma}\right) < 1\%. \quad (2.15)$$

Thus, if we want the probability of error be smaller than a certain value the following inequality must be satisfied :

$$f(D-1) > \alpha\Sigma, \quad (2.16)$$

where  $\alpha$  is a numerical constant independent of  $K$ ,  $D$ , and  $n$  (for 1% error  $\alpha \approx 2.576$ ). This translates into the following bound for the number of memories :

$$K < K^{\max} = \frac{1}{\alpha^2(2n-3)!!} D^{n-1}. \quad (2.17)$$

Thus, as long as the number of memories is smaller than  $K^{\max}$ , the network initialized in one of the memories remains there and the dynamics does not flow away from it. It turns out, this is precisely the point when associative memory recall breaks. If the number of memories is smaller than  $K^{\max}$  our network works as intended. Once  $K$  exceeds  $K^{\max}$ , reliable recall breaks. This does not mean that the network becomes useless in that regime. In fact, it instead becomes a generative model. We will discuss this aspect later.

### What have we learned so far?

- The number of memories  $K$  is upper bounded.
- The Memory storage capacity heavily depends on the shape of the energy function  $F(\cdot)$  and the shape of the activation function  $f(\cdot)$ .
- The sharper the energy peaks around memories – the larger the memory storage capacity.

## 2.2 Limiting Cases

It is instructive to study a few limiting cases of the general family Eq. (2.1). Each of these models are frequently studied in the literature and have distinct properties.

**The Hopfield Model**  $n = 2$ . The simplest, and the most popular, example of the Dense Associative Memory is the Hopfield model. One can obtain it from the general form Eq. (2.1)

choosing the function as  $F(\cdot) = \frac{1}{2}(\cdot)^2$ . The energy function can be written as

$$E = -\frac{1}{2} \sum_{\mu=1}^K \left( \sum_{i=1}^D \xi_i^\mu \sigma_i \right)^2 = -\frac{1}{2} \sum_{i,j=1}^D \sigma_i T_{ij} \sigma_j, \quad \text{where} \quad T_{ij} = \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu. \quad (2.18)$$

In this case, according to the general result Eq. (2.17), the memory storage capacity scales linearly with the size of the network:

$$K^{\max} \sim D. \quad (2.19)$$

This is the famous  $K^{\max} \approx 0.14D$  scaling law from the Hopfield's 1982 paper [7], derived by [15] using tools from statistical mechanics. While this model is appealing from the perspective of mathematical elegance and simplicity, this scaling law presents a major practical limitation. In the end, the hallmark of modern AI applications is the ability to store and process large amounts of information, a property severely limited by this scaling law.

**DenseAM with  $n = 3$ .** Fortunately, this problem disappears for a more rapidly peaking energy function (obtained via an alternative activation function). For  $F(\cdot) = \frac{1}{3}(\cdot)^3$ , for example, the energy is given by

$$E = -\frac{1}{3} \sum_{\mu=1}^K \left( \sum_{i=1}^D \xi_i^\mu \sigma_i \right)^3 = -\frac{1}{3} \sum_{i,j,k=1}^D T_{ijk} \sigma_i \sigma_j \sigma_k, \quad \text{where} \quad T_{ijk} = \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu \xi_k^\mu, \quad (2.20)$$

and the memory storage capacity scales as:

$$K^{\max} \sim D^2, \quad (2.21)$$

which is significantly faster than linearly.

**DenseAM with  $F(\cdot) = \exp(\cdot)$ .** It turns out that one can even achieve the exponentially large memory storage capacity. For exponential function  $F(\cdot)$  [27; 28], the number of memories that this DenseAM can store and retrieve scale as:

$$K^{\max} \sim 2^{\frac{D}{2}}, \quad (2.22)$$

which is more than sufficient for storing any practically relevant amount of information. Note, this number is the square root of the total number of binary states of the network. Despite its huge memory storage capacity, this model retains strong error correcting capabilities and has large size basins of attraction around each stored memories.

## 2.3 General Dense Associative Memory with Binary State Variables

Although simple models represented by Eq. (2.1) illustrate the computational capacities of Dense Associative Memories, more general energy functions are also frequently studied. For binary

DenseAM models, the general form of the energy function is given by

$$E = -Q \left[ \sum_{\mu=1}^K F \left( S[\xi^\mu, \sigma] \right) \right], \quad (2.23)$$

where the function  $F(\cdot)$  is a rapidly growing separation function (e.g., power  $F(\cdot) = (\cdot)^n$  or exponent),  $S[\mathbf{x}, \mathbf{x}']$  is a similarity function (e.g., a dot product or a Euclidean distance), and  $Q$  is a scalar monotone function (e.g., linear or logarithm). There are many possible combinations of various functions  $F(\cdot)$ ,  $S(\cdot, \cdot)$ , and  $Q(\cdot)$  that lead to different models from the DenseAM family [21; 27; 29; 30; 31; 32]. We will discuss the relationship between these binary models and DenseAMs with continuous states in the next Chapter.

### Notebook 2.1: Storage and recovery of memories in DenseAM

In this notebook, we offer the reader the possibility to experience storage and retrieval of patterns in DenseAM models. A set of simple Pokemon images can be embedded in the memory pool of the model. The model can then be queried by a corrupted version of a memory. The dynamical trajectory of the recall process retrieves the desired memory. By varying parameter  $n$  the reader can experience both successful recovery of the memories and memory failures, when the recovered image does not correspond to the desired memory. All these numerical results are to illustrate the general theory discussed in this chapter.

Checkout the notebook as a [blog post](#), a [colab notebook](#) or as a [raw .ipynb](#) file.



**Minimize  
energy**

How many patterns  
can we store?



## Chapter 3

# General Dense Associative Memory

In the previous Chapter, we introduced Dense Associative Memories with discrete state vectors. While mathematically aesthetic and simple to analyze, such models do not allow backpropagation training — due to their discrete nature. It turns out, most of the desired properties of Dense Associative Memories with discrete states are inherited by models with continuous variables. Moreover, the discrete state models can be derived as limiting cases of the continuous models.

There are two other limitations of the models represented by Eq. (2.1). First, they do not have the hierarchical structure of representations, a crucial aspect which limits their ability to handle complex patterns from real-world datasets. Second, they have a rigid energy function — although the energy depends on the learnable parameters  $\xi_i^\mu$  — its specific form may constrain the types of patterns and relationships the network can model.

In this Chapter, we introduce “building blocks” of Dense Associative Memories. Specifically, we develop a *modular energy* perspective where the energy of any model from this family can be decomposed into standardized components: **neuron layers** that encode dynamic variables and **hypersynapses** that encode their interactions. *The total energy of the system is the sum of the individual component energies subtracted by the energies of all the layers and all the interactions between those layers.* This framework of energy-based building blocks for memory not only clarifies how existing methods relate to each other, but also provides a systematic language for designing new architectures. We demonstrate the flexibility of this abstraction by showing how it helps to formulate all of the known models from this family, including Hierarchical Associative Memories [33], Energy Transformers [34], neuron-astrocyte networks [35], and many others.

We refer to this generalized abstraction of Energy-based AMs as **HAMUX** [36] after the software library that introduced it (here, HAMUX stands for “**H**ierarchical **A**ssociative **M**emory **U**ser **eX**perience”). We emphasize, however, that the abstraction is more fundamental than its specific software implementation.

### 3.1 Building Blocks of AMs with Modular Energies

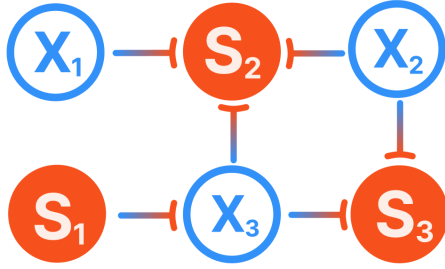
HAMUX builds deep AMs by summing the modular energies of *neuron layers* and *hypersynapses*.

A **neuron layer** captures a non-linearity in the network (e.g., ReLU, sigmoid, tanh, softmax, layernorm, etc.). We call these non-linearities *activations*, and their inputs or *pre-activations* serve as the dynamic variables of the system. For example, a neuron layer can capture the

## Associative Memory

A hypergraph of **neurons** communicating via **synapses**

**One** total energy  
**Local** computation  
**Guaranteed** convergence



$$E_{\text{total}} = \sum_{\text{neurons}} E + \sum_{\text{synapses}} E$$

$$\frac{d\mathbf{x}_\ell}{dt} = -\frac{\partial E_{\text{total}}}{\partial \hat{\mathbf{x}}_\ell}$$

### Neuron Layer

Simplify non-linear dynamics using **Lagrangians**

Internal **States**  $\mathbf{X}$   $\xrightarrow{\mathcal{L}}$   $\hat{\mathbf{X}}$  **Non-linear Activations**

Legendre Transform of the Lagrangian defines both **activations** and layer's **energy**

$$\hat{\mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \quad \Bigg| \quad E = \langle \mathbf{x}, \hat{\mathbf{x}} \rangle - \mathcal{L}$$

Dynamic **states** minimize energy of activations

$$\frac{d\mathbf{x}}{dt} = -\frac{\partial E}{\partial \hat{\mathbf{x}}}$$

### Hypersynapse

Learn to align activations of connected neurons



Low energy means aligned activations.  
 Minimizing energy maximizes "similarity"

$$E = -\text{sim}(\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots; \Xi)$$

**Figure 3.1:** HAMUX hypergraph diagrams are a graphical depiction of an AM whose total energy is the sum of the **neuron layer** (node) and **hypersynapse** (hyperedge) energies. Inference is done recurrently, modeled by a system of differential equations where each neuron layer's hidden state updates to minimize the total energy. When all non-linearities are captured in the dynamic neurons, inference becomes a local computation that avoids differentiating through non-linearities.

computation  $\hat{\mathbf{x}} = \text{ReLU}(\mathbf{x})$ , which has activations  $\hat{\mathbf{x}}$  and pre-activations  $\mathbf{x}$  which serve as the dynamic *internal state* for this neuron layer. Structurally, neuron layers are the *nodes* of our energy-based computation graph.

A **hypersynapse** is a parameterized energy function that captures how similar or *aligned* the activations of its connected neuron layers are. For example, a simple hypersynapse may take the form  $E_S(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \Xi) = \hat{\mathbf{x}}^\top \Xi \hat{\mathbf{y}}$ , where  $\Xi$  is a synaptic weight matrix. The gradient of  $E_S$  w.r.t.  $\hat{\mathbf{x}}$  or  $\hat{\mathbf{y}}$  looks like a Dense linear transformation, though more complex synaptic energies can be chosen to look like Conv, Pooling, or even Attention layers. Hypersynapses define the interactions between neurons and are the *hyperedges* of our energy-based computation graph.

For a system of  $L$  neuron layers and  $S$  hypersynapses, the total energy of the system is

$$E_{\text{total}} = \sum_{\ell=1}^L E_{\ell}^{\text{neuron}} + \sum_{s=1}^S E_s^{\text{synapse}}. \quad (3.1)$$

The total energy is structured such that the activations of a neuron layer affect only connected hypersynapses and itself. Let  $\hat{\mathbf{x}}_{\ell}$  and  $\mathbf{x}_{\ell}$  represent the activations and internal states of neuron layer  $\ell$ , and let  $\mathbf{N}(\ell)$  represent the set of hypersynapses that connect to neuron layer  $\ell$ . The following update rule describes how neuron internal states  $\mathbf{x}_{\ell}$  minimize the total energy using only local signals

$$\tau_{\ell} \frac{d\mathbf{x}_{\ell}}{dt} = -\frac{\partial E_{\text{total}}}{\partial \hat{\mathbf{x}}_{\ell}} = -\left( \sum_{s \in \mathbf{N}(\ell)} \frac{\partial E_s^{\text{synapse}}}{\partial \hat{\mathbf{x}}_{\ell}} \right) - \frac{\partial E_{\ell}^{\text{neuron}}}{\partial \hat{\mathbf{x}}_{\ell}} = \mathcal{I}_{x_{\ell}} - \mathbf{x}_{\ell}, \quad (3.2)$$

where  $\mathcal{I}_{x_{\ell}} := -\sum_{s \in \mathbf{N}(\ell)} \nabla_{\hat{\mathbf{x}}_{\ell}} E_s^{\text{synapse}}$  is the *total synaptic input current* to neuron layer  $\ell$ , which is fundamentally local and serves to minimize the energy of connected hypersynapses. See sections (3.2) and (3.3) to understand the above equation in more detail. The time constant for neurons in layer  $i$  is denoted by  $\tau_{\ell}$ . When the activations  $\hat{\mathbf{x}}_{\ell}$  are bounded, the above system is guaranteed to converge for any choice of hypersynapse energies.

## 3.2 Dynamical Neurons and their Lagrangians

A *neuron layer* or *node* is a fancy term to describe the dynamic variables in AM. Each neuron layer has an *internal state*  $\mathbf{x}$  which evolves over time and an *activation*  $\hat{\mathbf{x}}$  that forwards a signal to the rest of the network. Think of neurons like the activation functions of standard neural networks, where  $\mathbf{x}$  are the “pre-activations” and  $\hat{\mathbf{x}}$  are the outputs e.g.,  $\hat{\mathbf{x}} = \text{ReLU}(\mathbf{x})$ .

In order to define neuron’s layer energy, AMs employ two mathematical tools from physics: *convex Lagrangian functions* and the *Legendre transform*. For each neuron layer, we define a convex, scalar-valued Lagrangian  $\mathcal{L}_x(\mathbf{x})$ . The Legendre transform  $\mathcal{T}$  of this Lagrangian produces the dual variable  $\hat{\mathbf{x}}$  (our activations) and the dual energy  $E_x(\hat{\mathbf{x}})$  (our new energy) as in:

$$\begin{aligned} \hat{\mathbf{x}} &= \nabla \mathcal{L}_x(\mathbf{x}) \quad (\text{activation function}) \\ E_x(\hat{\mathbf{x}}) &= \mathcal{T}(\mathcal{L}_x) = \langle \mathbf{x}, \hat{\mathbf{x}} \rangle - \mathcal{L}_x(\mathbf{x}) \quad (\text{dual energy}) \end{aligned} \quad (3.3)$$

where  $\langle \cdot, \cdot \rangle$  is the element-wise inner product. Because  $\mathcal{L}_x$  is convex, the Jacobian of the activations  $\frac{\partial \hat{\mathbf{x}}}{\partial \mathbf{x}} = \nabla^2 \mathcal{L}_x(\mathbf{x})$  (i.e., the Hessian of the Lagrangian) is positive definite. This important point is summarized in Fig. (3.1).

The dual energy  $E_x(\hat{\mathbf{x}})$  has another nice property: *its gradient equals the hidden states*. Thus, when we minimize the energy of our neurons (in the absence of any other signal), we observe exponential decay. This is nice to keep the dynamic behavior of our system bounded and well-behaved, especially for very large values of  $\mathbf{x}$ .

$$\frac{d\mathbf{x}}{dt} = -\nabla_{\hat{\mathbf{x}}} E_x(\hat{\mathbf{x}}) = -\mathbf{x}. \quad (3.4)$$

**Summary** For non-physicists, the terminology used in this section can be daunting. The key insight is simple: a neuron layer is just a convex function  $\mathcal{L}_x$  (the Lagrangian) applied to an internal state  $\mathbf{x}$ . The Legendre transform of this Lagrangian then automatically provides two things: (1) the activation function  $\hat{\mathbf{x}} = \nabla \mathcal{L}_x(\mathbf{x})$ , and (2) the dual energy representation  $E_x(\hat{\mathbf{x}})$ . This mathematical machinery abstracts away some of the complexity of non-linearities and gives us a simpler system to work with.

**Proof: Energy gradient equals hidden states**

Show that  $\frac{\partial E_x(\hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}} = \mathbf{x}$ .

$$\begin{aligned} \frac{\partial E_x(\hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}} &= \frac{\partial}{\partial \hat{\mathbf{x}}} (\langle \mathbf{x}, \hat{\mathbf{x}} \rangle - \mathcal{L}_x(\mathbf{x})) \\ &= \mathbf{x} + \hat{\mathbf{x}} \frac{\partial \mathbf{x}}{\partial \hat{\mathbf{x}}} - \frac{\partial \mathcal{L}_x(\mathbf{x})}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \hat{\mathbf{x}}} \\ &= \mathbf{x} + \hat{\mathbf{x}} \frac{\partial \mathbf{x}}{\partial \hat{\mathbf{x}}} - \hat{\mathbf{x}} \frac{\partial \mathbf{x}}{\partial \hat{\mathbf{x}}} \\ &= \mathbf{x} \end{aligned}$$

### 3.3 Hypersynapses

The activations of one neuron layer are sent to other neurons via communication channels called *hypersynapses*. At its most general, a hypersynapse is a scalar valued energy function defined on top of the activations of connected neuron layers. For example, a hypersynapse connecting neuron layers  $\mathbf{X}$  and  $\mathbf{Y}$  has an *interaction energy*  $E_{xy}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \Xi)$ , where  $\Xi$  represents the *synaptic weights* or learnable parameters.  $E_{xy}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \Xi)$  encodes the desired relationship between activations  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{y}}$ : when this energy is low, the activations satisfy the relationship encoded by the synaptic weights  $\Xi$ . During energy minimization, the system adjusts the activations to reduce all energy terms, which means synapses effectively “pull” the connected neuron layers toward configurations encoded in the parameters that minimize their interaction energy.

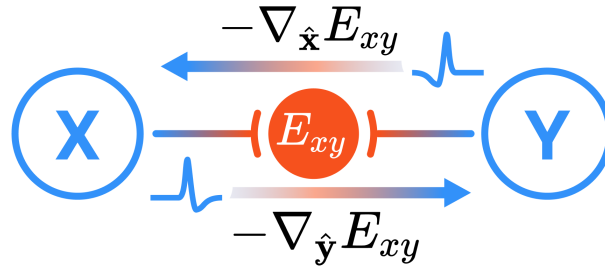
Hypersynapses in the HAMUX framework differ from biological synapses in two fundamental ways:

1. **Hypersynapses can connect any number of layers simultaneously**, while biological synapses connect only two neurons. This officially makes hypersynapses “hyperedges” in graph theory terms.
2. **Hypersynapses are undirected**, meaning that all connected layers influence each other bidirectionally during energy minimization. Meanwhile, biological synapses are unidirectional, meaning signal flows from a presynaptic to postsynaptic neuron.

Because of these differences, we choose the distinct term “hypersynapses” to distinguish them from biological synapses.

### Hypersynapses are undirected edges

Connecting two neurons sends signal both ways



**Figure 3.2:** Hypersynapses are represented as undirected (hyper)edges in a hyper-graph. Shown is an example pairwise synapse, which is a single energy function  $E_{xy}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \Xi)$  defined on the activations  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{y}}$  from connected nodes, which necessarily propagate signal to both connected nodes. Here, *signal* is defined as the negative gradient of the interaction energy w.r.t. the connected layer’s activations (e.g., layer X receives signal  $\mathcal{I}_x = -\nabla_{\hat{\mathbf{x}}} E_{xy}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \Xi)$ ). This is in contrast to biological synapses which are directional and only propagate signal in one direction from layer X to Y, needing a separate synapse to bring information back from Y to X.

#### Hypersynapse notation conventions

For synapses connecting multiple layers, we subscript with the identifiers of all connected layers. For example:

- $E_{xy}$  — synapse connecting layers X and Y
- $E_{xyz}$  — synapse connecting layers X, Y, and Z.
- $E_{xyz\dots}$  — synapses connecting more than three layers are possible, but rare.

However, synapses can also connect a layer to itself (self-connections). To avoid confusion with neuron layer energy  $E_x$ , we use curly brackets for synaptic self-connections. For example,  $E_{\{x\}}$  represents the interaction energy of a synapse that connects layer X to itself.

Because almost every interaction energy is parameterized in some way, we generally omit  $\Xi$  from the notation in subsequent sections when it’s not central to the discussion

The undirected nature of hypersynapses fundamentally distinguishes AM from traditional neural networks. Whereas feed-forward networks follow a directed computational graph with clear input-to-output flow, AMs have no inherent concept of “forward” or “backward” directions. All connected layers influence each other bidirectionally during energy minimization, with information propagating from deeper layers to shallower layers as readily as the other way around. See Fig. (3.1) for a visual illustration.

Unlike the neuron layer’s energies, the interaction energies of the hypersynapses are completely unconstrained: *any function* that takes activations as input and returns a scalar is admissible and will have well-behaved dynamics<sup>1</sup>. The interaction energy of a synapse may choose to introduce

<sup>1</sup>Some energies could be more meaningful than the others.

its own non-linearities beyond those handled by the neuron layers. When this occurs, the energy minimization dynamics must compute gradients through these “synaptic non-linearities”, unlike the case where all non-linearities are abstracted into the neuron layer Lagrangians.

### 3.4 Energy Descent Dynamics

The central result is that dynamical equations Eq. (3.2) decrease the global energy of the network Eq. (3.1). In order to demonstrate this, consider the total time derivative of the energy

$$\frac{dE_{\text{total}}}{dt} = \sum_{i=1}^L \frac{\partial E_{\text{total}}}{\partial \hat{\mathbf{x}}_i} \frac{\partial \hat{\mathbf{x}}_i}{\partial \mathbf{x}_i} \frac{d\mathbf{x}_i}{dt} = - \sum_{i=1}^L \tau_i \frac{d\mathbf{x}_i}{dt} \frac{\partial^2 \mathcal{L}_x}{\partial \mathbf{x}_i \partial \mathbf{x}_i} \frac{d\mathbf{x}_i}{dt} \leq 0, \quad (3.5)$$

where we expressed the partial of the energy w.r.t. the activations through the velocity of the neuron’s internal states Eq. (3.2). The Hessian matrix  $\frac{\partial^2 \mathcal{L}_x}{\partial \mathbf{x}_i \partial \mathbf{x}_i}$  has the size number of neurons in layer  $i$  multiplies by the number of neurons in layer  $i$ . As long as this matrix is positive semi-definite, a property resulting from the convexity of the Lagrangian, the total energy of the network is guaranteed to either decrease or stay constant — increase of the energy is not allowed.

Additionally, if the energy of the network is bounded from below, the dynamics in Eq. (3.2) are guaranteed to lead the trajectories to fixed manifolds corresponding to local minima of the energy. If the fixed manifolds have zero-dimension, i.e., they are fixed point attractors, the velocity field will vanish once the network arrives at the local minimum. This corresponds to Hessians being strictly positive definite. Alternatively, if the Lagrangians have zero modes, resulting in existence of zero eigenvalues of the Hessian matrices, the network may converge to the fixed manifolds, but the velocity fields may stay non-zero, while the network’s state moves along that manifold.

### 3.5 Implementing AMs

We have established how the computational graph is built and the rules for how neuron layers and hypersynapses are constructed. We now discuss how the above mathematical framework can be used to recreate some of the commonly used AM models.

#### Exercise 3.1: Designing the energy for a custom DenseAM

**Problem** Consider a DenseAM model consisting of  $D$  neurons with the following activation function  $\hat{x}_i = \tanh(\beta x_i)$ . Design the synaptic energy and the global energy to recreate DenseAM with discrete variables discussed in Eqs. (2.1) and (2.3), in the limit  $\beta \rightarrow \infty$ .

#### Solution

First, define the Lagrangian for this network so that its partial gives the desired activation

$$\mathcal{L} = \frac{1}{\beta} \sum_{i=1}^D \log \left( \cosh(\beta x_i) \right), \quad \text{resulting in} \quad \hat{x}_i = \tanh(\beta x_i) \quad (3.6)$$

The synapse connects neuron layer to itself and its synaptic energy is given by

$$E^{\text{synapse}} = - \sum_{\mu=1}^K F\left(\sum_{j=1}^D \xi_j^\mu \hat{x}_j\right) \quad (3.7)$$

The total energy of the network is

$$E^{\text{total}} = E^{\text{neuron}} + E^{\text{synapse}} = \left[ \sum_{i=1}^D \hat{x}_i x_i - \mathcal{L} \right] - \sum_{\mu=1}^K F\left(\sum_{j=1}^D \xi_j^\mu \hat{x}_j\right) \quad (3.8)$$

The dynamical update equation Eq. (3.2) is given by

$$\tau \frac{dx_i}{dt} = \sum_{\mu=1}^K \xi_i^\mu f\left(\sum_{j=1}^D \xi_j^\mu \hat{x}_j\right) - x_i \quad (3.9)$$

Now, let's discretize time. Set  $\tau = 1$  and write the above equation in finite differences ( $dt = 1$ ). The result is

$$\frac{x_i^{t+1} - x_i^t}{dt} = \sum_{\mu=1}^K \xi_i^\mu f\left(\sum_{j=1}^D \xi_j^\mu \hat{x}_j^t\right) - x_i^t \quad (3.10)$$

which leads to

$$x_i^{t+1} = \sum_{\mu=1}^K \xi_i^\mu f\left(\sum_{j=1}^D \xi_j^\mu \hat{x}_j^t\right) \quad (3.11)$$

Finally, express everything through the activations  $\hat{x}_i$  and take the limit  $\beta \rightarrow \infty$ . In this limit  $\hat{x}_i = \sigma_i = \text{Sign}(x_i)$  and the energy of the layer vanishes, resulting in the total energy

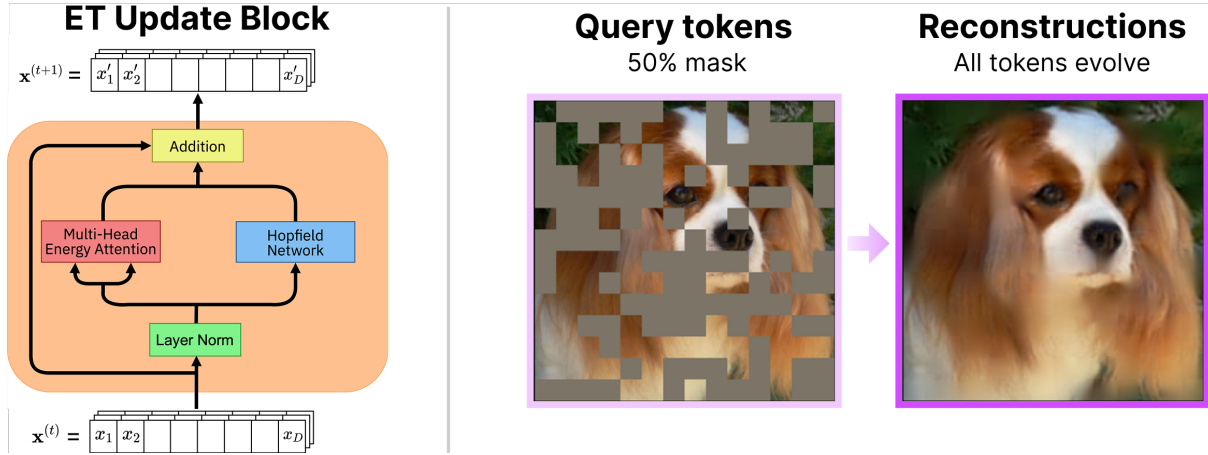
$$E^{\text{total}} = - \sum_{\mu=1}^K F\left(\sum_{j=1}^D \xi_j^\mu \sigma_j^t\right) \quad (3.12)$$

The discrete update equation can be obtained by acting with the  $\text{Sign}(\cdot)$  function on both sides of Eq. (3.11), resulting in Eq. (2.3).

### 3.5.1 Energy Transformer Block

We now explain how the techniques developed above can be used for building the Energy Transformer (ET) architecture [34]. For clarity of presentation, we use language associated with the image domain, although this architecture can also be used for language or graphs with minimal modifications.

The overall pipeline is similar to the Vision Transformer networks (ViTs) [37]. An input image is split into non-overlapping patches. After passing these patches through the encoder and adding the positional information, the semantic content of each patch and its position is encoded in the token  $x_{iA}$ . In the following the indices  $i, j, k = 1, \dots, D$  are used to denote the token vector's elements,



**Figure 3.3:** Inside the ET block. The input tokens  $\mathbf{x}^{(t)}$  passes through a sequence of operations and gets updated to produce the output tokens  $\mathbf{x}^{(t+1)}$ . The operations inside the ET block are carefully engineered so that the entire network has a global energy function, which decreases with time and is bounded from below. In contrast to conventional transformers, the ET-based analogs of the attention module and the feed-forward MLP module are applied in parallel as opposed to consecutively. **Right:** The ET block recurrently minimizes the energy of a corrupted image represented by a collection of tokens, where 50% of the tokens are occluded. Shown is an image not seen when training the ET block.

indices  $A, B, C = 1, \dots, N$  are used to enumerate the patches and their corresponding tokens. It is helpful to think about each image patch as a physical particle, which has a complicated internal state described by a  $D$ -dimensional vector  $\mathbf{x}_A$ . This internal state describes the identity of the particle (representing the pixels of each patch), and the particle’s positional embedding (the patch’s location within the image). The ET block is described by a continuous time differential equation, which describes interactions between these particles. Initially, at  $t = 1$  the network is given a set containing two groups of particles corresponding to open and masked patches. The “open” particles know their identity and location in the image. The “masked” particles only know where in the image they are located, but are not provided the information about what image patch they represent. The goal of ET’s non-linear dynamics is to allow the masked particles to find an identity consistent with their locations and the identities of open particles. This dynamical evolution is designed so that it minimizes a global energy function. The identities of the masked particles are considered to be revealed when the dynamical trajectory reaches the fixed point. Thus, the central question is: how can we design the energy function that accurately captures the task that the Energy Transformer needs to solve?

The masked particles’ search for identity is guided by two pieces of information: identities of the open particles, and the general knowledge about what patches are in principle possible in the space of all possible images. These two pieces of information are described by two contributions to the ET’s energy function: the energy based attention and the Hopfield Network. Below we define each element of the ET block in the order they appear in Fig. (3.3).

### Layer-Norm

Each token, or a particle, is represented by a vector  $\mathbf{x} \in \mathbb{R}^D$ . At the same time, most of the operations inside the ET block are defined using a layer-normalized token representation:

$$\hat{x}_i = \gamma \frac{x_i - \bar{x}}{\sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon}} + \delta_i, \quad \text{where} \quad \bar{x} = \frac{1}{D} \sum_{k=1}^D x_k \quad (3.13)$$

The scalar  $\gamma$  and the vector elements  $\delta_i$  are learnable parameters,  $\varepsilon$  is a small regularization constant. Following the general recipe of HAMUX, this operation can be viewed as an activation function for the neural layer and can be derived as a partial derivative of the Lagrangian function:

$$\mathcal{L}(\mathbf{x}) = D\gamma \sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon} + \sum_j \delta_j x_j, \quad \text{so that} \quad \hat{x}_i = \frac{\partial \mathcal{L}(\mathbf{x})}{\partial x_i} \quad (3.14)$$

See [30; 38; 33] for the discussion of this property.

### Multi-Head Energy Attention

The first contribution to the ET's energy function is responsible for exchanging information between the particles (tokens). Similarly to the conventional attention mechanism, each token generates a pair of queries and keys (ET does not have a separate value matrix; instead the value matrix is a function of keys and queries). The goal of the energy based attention is to evolve the tokens in such a way that the keys of the open patches are aligned with the queries of the masked patches in the internal space of the attention operation. Below we use index  $\alpha = 1, \dots, Y$  to denote elements of this internal space, and index  $h = 1, \dots, H$  to denote different heads of this operation. With these notations the energy-based attention operation is described by the following energy function:

$$E^{\text{ATT}} = -\frac{1}{\beta} \sum_{h=1}^H \sum_{C=1}^N \log \left( \sum_{B \neq C} \exp(\beta A_{hBC}) \right) \quad (3.15)$$

where the attention matrix  $A_{hBC}$  is computed from query and key tensors as follows:

$$\begin{aligned} A_{hBC} &= \sum_{\alpha} K_{\alpha hB} Q_{\alpha hC}, & \mathbf{A} &\in \mathbb{R}^{H \times N \times N} \\ K_{\alpha hB} &= \sum_j W_{\alpha h j}^K \hat{x}_{jB}, & \mathbf{K} &\in \mathbb{R}^{Y \times H \times N} \\ Q_{\alpha hC} &= \sum_j W_{\alpha h j}^Q \hat{x}_{jC}, & \mathbf{Q} &\in \mathbb{R}^{Y \times H \times N} \end{aligned} \quad (3.16)$$

and the tensors  $\mathbf{W}^K \in \mathbb{R}^{Y \times H \times D}$  and  $\mathbf{W}^Q \in \mathbb{R}^{Y \times H \times D}$  are learnable parameters. From the perspective of HAMUX, Eq. (3.15) is the energy of the synapse, which mixes layers of neurons or tokens.

From the computational perspective each patch generates two representations: query (given the

position of the patch and its current content, where in the image should it look for the prompts on how to evolve in time?), and key (given the current content of the patch and its position, what should be the contents of the patches that attend to it?). The log-sum energy function (3.15) is minimal when for every patch in the image its queries are aligned with the keys of a small number of other patches connected by the attention map. Different heads (index  $h$ ) contribute to the energy additively.

### Hopfield Network Module

The next step of the ET block, which we call the Hopfield Network (HN), is responsible for ensuring that the token representations are consistent with what one expects to see in realistic images. The energy of this sub-block is defined as:

$$E^{\text{HN}} = - \sum_{B=1}^N \sum_{\mu=1}^K G\left(\sum_{j=1}^D \xi_{\mu j} \hat{x}_{jB}\right), \quad \boldsymbol{\xi} \in \mathbb{R}^{K \times D} \quad (3.17)$$

where  $\xi_{\mu j}$  is a set of learnable weights (memories in the Hopfield Network), and  $G(\cdot)$  is an integral of the activation function  $r(\cdot)$ , so that  $G(\cdot)' = r(\cdot)$ . This formula is identical to the energy Eq. (3.7). Depending on the choice of the activation function this step can be viewed either as a classical continuous Hopfield Network [14] if the activation function grows slowly (e.g.,  $r(\cdot) = \text{ReLU}$ ), or as a Dense Associative Memory [21; 30] if the activation function is sharply peaked around the memories (e.g.,  $r(\cdot) = \text{power}$  or  $\text{softmax}$ ). The HN sub-block is analogous to the feed-forward MLP step in the conventional transformer block but requires that the weights of the projection from the token space to the hidden neurons' space to be the same (transposed matrix) as the weights of the subsequent projection from the hidden space to the token space. Thus, the HN module here is an MLP with shared weights that is *applied recurrently*. The energy contribution of this block is low when the token representations are aligned with some rows of the matrix  $\xi$ , which represent memories, and high otherwise.

### Dynamics of Token Updates

The inference pass of the ET network is described by the continuous time differential equation, which minimizes the sum of the two energies described above. The whole ET network contains of layers of tokens coupled through two types of synapses, attention synapse and Hopfield Network synapse, so that

$$\begin{aligned} E_{\text{total}} &= E^{\text{neuron}} + \sum_{\alpha=1}^2 E_{\alpha}^{\text{synapse}} \\ &= \left[ \sum_{A=1}^N \sum_{i=1}^D x_{iA} \hat{x}_{iA} - \sum_{A=1}^N \mathcal{L}(\mathbf{x}_A) \right] + E^{\text{ATT}} + E^{\text{HN}} \\ &\approx E^{\text{ATT}} + E^{\text{HN}} + O(\varepsilon) \end{aligned} \quad (3.18)$$

We work in the regime when the parameter  $\varepsilon$  in the definition of the layer-norm Lagrangian is small — it only serves as a regularization to prevent the division by zero. In this limit, neuron

layer energy vanishes, and the total HAMUX energy is the sum of  $E^{\text{ATT}}$  and  $E^{\text{HN}}$

$$\tau \frac{dx_{iA}}{dt} = -\frac{\partial E_{\text{total}}}{\partial \hat{x}_{iA}}, \quad \text{where} \quad E_{\text{total}} = E^{\text{ATT}} + E^{\text{HN}} \quad (3.19)$$

Here  $x_{iA}$  is the token representation (input and output from the ET block), and  $\hat{x}_{iA}$  is its layer-normalized version. The first energy is low when each patch's queries are aligned with the keys of its neighbors. The second energy is low when each patch has content consistent with the general expectations about what an image patch should look like (memory slots of the matrix  $\xi$ ). The dynamical system, represented by Eq. (3.19), finds a trade-off between these two desirable properties of each token's representation. For numerical evaluations, Eq. (3.19) is discretized in time.

To demonstrate that the dynamical system (3.19) minimizes the energy, consider the temporal derivative

$$\frac{dE_{\text{total}}}{dt} = \sum_{i,j,A} \frac{\partial E_{\text{total}}}{\partial \hat{x}_{iA}} \frac{\partial \hat{x}_{iA}}{\partial x_{jA}} \frac{dx_{jA}}{dt} = -\frac{1}{\tau} \sum_{i,j,A} \frac{\partial E_{\text{total}}}{\partial \hat{x}_{iA}} M_{ij}^A \frac{\partial E_{\text{total}}}{\partial \hat{x}_{jA}} \leq 0 \quad (3.20)$$

The last inequality sign holds if the symmetric part of the matrix

$$M_{ij}^A = \frac{\partial \hat{x}_{iA}}{\partial x_{jA}} = \frac{\partial^2 \mathcal{L}}{\partial x_{iA} \partial x_{jA}} \quad (3.21)$$

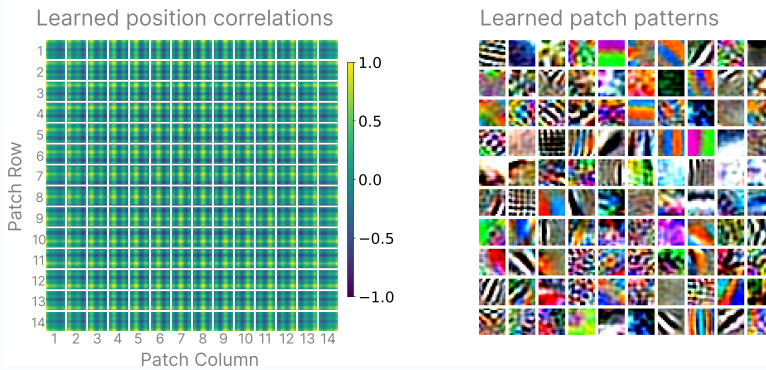
is positive semi-definite (for each value of index  $A$ ). The Lagrangian (3.14) satisfies this condition.

### Notebook 3.1: Energy Transformer

In this notebook, we offer the reader the possibility to build the ET block in code following the general rules of HAMUX. We have pre-trained this network on ImageNet and loaded the weights of the model, so that the reader can quickly play with parameters and visualize energy decent dynamics and learned representations at inference time. All these numerical results are to illustrate the general theory discussed in this chapter.

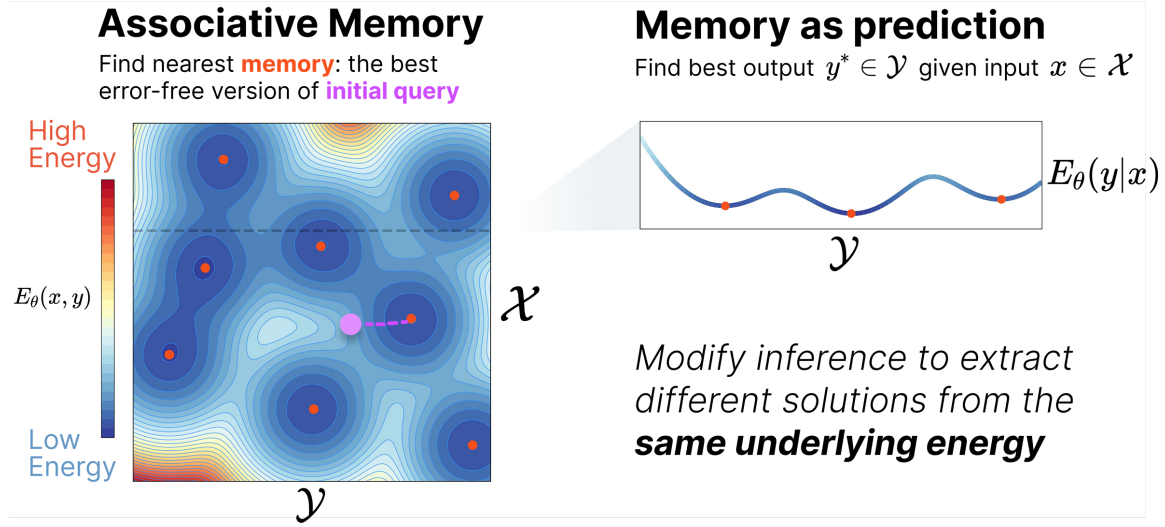
Checkout the notebook as a [blog post](#), a [colab notebook](#) or as a [raw .ipynb file](#).

#### Build and interpret the trained Energy Transformer



### 3.6 Bridging Energy Minimization and Feedforward Prediction

Although associative memories fundamentally differ from feedforward networks, they can both be used to solve the same tasks. Let  $\theta$  describe the network parameters. Traditional feedforward networks transform input tensors  $\mathbf{x} \in \mathcal{X}$  to output tensors  $\mathbf{y} \in \mathcal{Y}$  via  $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$ , where  $\mathbf{y}^* = f_\theta(\mathbf{x})$  represents the model’s prediction. In contrast, an AM builds an energy-based computational graph that maps input tensors  $\mathbf{z} \in \mathcal{Z}$  to a scalar energy value via  $E_\theta : \mathcal{Z} \mapsto \mathbb{R}$ .



**Figure 3.4: Associative Memory is fully compatible with traditional prediction tasks.** By fixing a subset of variables (input) and minimizing the energy with respect to the remaining variables, we can predict optimal values for the output variables. Left: a 2D energy landscape used for the general Associative Memory task of cleaning an input. Right: a slice through the total energy landscape represents the energy objective for a prediction task.

How can we use an AM for prediction? When input space  $\mathcal{X}$  and output space  $\mathcal{Y}$  are distinct (e.g., in classification or segmentation), the energy function takes both spaces as input:  $E_\theta : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$ . Given an input  $\mathbf{x} \in \mathcal{X}$  for which we want to know the best output  $\mathbf{y}^* \in \mathcal{Y}$  (as described in Eq. (3.22)), prediction or *inference* becomes a coordinate-wise constrained energy minimization problem where we fix one of the variables and minimize the energy with respect to the other.

$$\mathbf{y}^* = \arg \min_{\mathbf{y} \in \mathcal{Y}} E_\theta(\mathbf{x}, \mathbf{y}) \quad (3.22)$$

Sometimes the output space  $\mathcal{Y}$  of a feedforward network represents a noiseless, inpainted version of  $\mathcal{X}$ , as in masked-token prediction where  $\mathbf{x}^* = f_\theta(\mathbf{x}^{(0)})$ . In this case,  $f_\theta : \mathcal{X} \mapsto \mathcal{X}$  is doing a proxy of “error minimization” inside its computation graph (which likely has many residual connections), and the energy function is  $E_\theta : \mathcal{X} \mapsto \mathbb{R}$ .

The inference process of Eq. (3.22) is flexible and can be adapted for different prediction tasks. Let’s view the *optimization objective* of Eq. (3.22) as a higher order function  $\mathcal{F}$  applied to energy function  $E_\theta$ . In this way, we can describe different input-output mappings from the energy. For instance, Eq. (3.22) is a mapping  $\mathcal{F}_{x \rightarrow y}(E_\theta) : \mathcal{X} \mapsto \mathcal{Y}$  that performs a global search of  $E_\theta$  over possible  $\mathbf{y}$  given a clamped  $\mathbf{x}$ . We can define other useful inference objectives. For example,

we can invert the search to instead search over  $\mathbf{x}$  in the region of some initial guess  $\mathbf{x}^{(0)}$  given a clamped  $\mathbf{y}$ . This would be represented via  $\mathcal{F}_{(x,y) \rightarrow x}(E_{\theta}) : \mathcal{X} \times \mathcal{Y} \mapsto \mathcal{X}$ . Or we could jointly optimize both variables given only an initial guess for both, as in  $\mathcal{F}_{(x,y) \rightarrow (x,y)}(E_{\theta}) : \mathcal{X} \times \mathcal{Y} \mapsto \mathcal{X} \times \mathcal{Y}$ . The Hopfield Network generally considers this last scenario, but each inference process represents a different way to extract solutions from the same underlying energy landscape.

### Energy function vs. Loss function

The key idea in AMs is that every computation serves to optimize some objective. However, we choose to distinguish between two types of objectives: the *energy* function and the *loss* function. We say the energy function governs the dynamics of neuron states during *inference*, while the loss function governs the dynamics of model parameters during *training*. The primary difference is in whether the gradient is taken with respect to the states of the network, or the parameters.

## Chapter 4

# Failure of Memory and Generative AI

In 1977, psychologists Roger Brown and James Kulik described a famous experiment, in which respondents were asked to self-report the circumstances in which they found out about the highly surprising and consequential news of the President John F. Kennedy assassination [39]. Among many insightful findings from this study, peculiar responses have been recorded containing detailed, emotional, highly realistic, and convincing descriptions of learning about this news for the first time that were factually inaccurate. For instance, one of the respondents (person A) vividly describes how person B came down the stairs to the first floor of the house, while person A was focusing on work and told person A that she heard about the assassination on the news. This recollection is detailed enough to include specific phrases and portions of the conversation between persons A and B during this recalled event. Although both persons are well familiar with each other and their recollections are plausible, documented evidence suggests that person A and person B could not be present in the same location after the JFK’s assassination [40; 41].

This is an example of misremembering, a phenomenon that is general and can be observed in many other situations. For instance, during a crime investigation two eyewitnesses can give mutually contradictory accounts of what they saw. Sometimes, both accounts can be different from what has actually happened. These examples of misremembering highlight a failure mode of human memory in which multiple observed events (training data) blend together and form novel memories, which are different from any of the observed events (training data points). Misremembering leads to creation of novel hypothetical memories, which in certain aspects share a degree of similarity (correlated) to the training data, but are distinct from individual training instances. Thus, misremembering can lead to creativity.

In generative AI, *creativity* is a key objective. For instance, diffusion models, being trained on a sufficiently large set of training images, can generate genuinely novel photorealistic images of previously unseen events [42; 43; 44; 45]. A typical diffusion model training pipeline contains of two phases: the forward process when the noise is injected into the training samples, and the reverse process when a neural network is used to predict how much noise should be removed from the noisy sample with the goal to reconstruct the original uncorrupted training data point. At runtime, random noise is passed through the reverse process and converted into generated samples. The training pipeline of diffusion models can be conceptualized as the process of writing the training data into a memory network [46; 47]. By doing so, the information about training samples is written into the synaptic weights of the neural network that is used for denoising. The reverse process can be conceptualized as an attempt of memory recall — that memorized

information should be retrieved from the synaptic weights and turned into a generated sample. It is well established that the memory recall from the diffusion models can be successful; in that case the generated sample matches exactly at least one of the training samples. It can also be unsuccessful; in that case the generated sample will be novel and will not match any of the training samples. This is what is called the memorization-to-generalization transition in diffusion models [48; 49; 50], which occurs when the size of the training set is increased. Successful memory retrieval is typically viewed as a negative outcome in diffusion models and can often lead to privacy and copyright violations. Similarly, in LLMs training samples can often be extracted verbatim from the synaptic weights of the neural network, a property that has been a subject of intense discussions in the research community and public discourse [48]. At the same time, LLMs can also generate novel previously unseen responses. Importantly for us, in all these examples *creativity arises as a result of a failure of memory recall*.

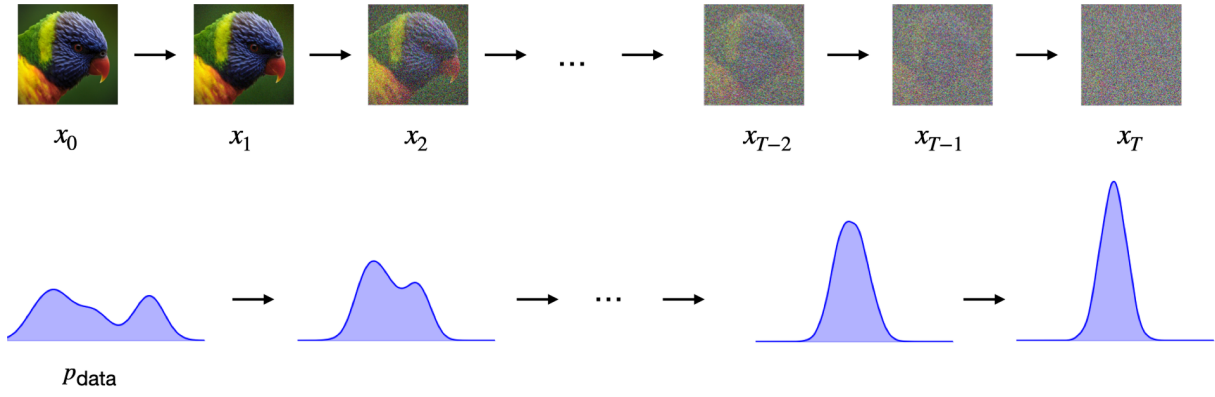
In Energy-based Associative Memory networks, memories are identified as local minima of an energy landscape, and the process of memory recall is conceptualized as a dynamical trajectory starting at a high energy state (corrupted memory) and leading to the best matching local minimum (recovered memory). Misremembering arises when the recovered memory (local minimum) of the energy function is different from any of the training data. These misremembered local minima are called *spurious states*, see the below Fig. (4.2) for their illustration. In AM literature, they are typically viewed as an obstacle to the faithful memory recall. For this reason, researchers in this field typically aim to either remove them entirely from the energy landscape, or mitigate their contribution to computation (e.g., by raising their energy) [20; 51].

In this Chapter, we discuss the emergence of *spurious states* in Dense Associative Memories, and the general relationship between these memory models and diffusion models, popular in generative AI. In previous chapters, AM models were studied in two situations. First, when the models have small memory capacity and are trained on a small amount of data, e.g., classical Hopfield Networks. Second, in situations when the models have large memory storage capacity, but still are trained on a small amount of data. The focus of this Chapter are settings in which the models are big (large information storage capabilities), and the amount of training data is even bigger (exceeds the critical memory storage capacity of the model). In this regime, DenseAMs turn into generative AI models.

## 4.1 Diffusion Models

Diffusion models have recently gained popularity, due to their flexibility and accuracy in modeling high-dimensional distributions for a variety of domains, including image generation [43; 45; 44], audio [52; 53; 54], video synthesis [55; 56; 57; 58], and other scientific applications. However, these powerful and flexible models pose great challenges related to privacy and security, as concerns grow about their tendency to generate their training data [49; 59; 48]. Such matters consequently emphasize the need for further understanding of memorization and generalization behaviors in diffusion models.

There are two fundamental processes which govern the aspects of diffusion models. Firstly, the *forward process* typically described by the following Itô Stochastic Differential Equation (SDE)



**Figure 4.1:** A general illustration of diffusion models. Addition of noise transforms the complex data distribution into a simple distribution — an isotropic Gaussian. The reverse process removes the noise and transforms a noise sample into a sample from the data distribution.

[45]:

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)dt + g(t)d\mathbf{w}_t, \quad (4.1)$$

transforms the given data distribution ( $\mathbf{x}_0 = \mathbf{y}$ ) into a simpler distribution<sup>1</sup>, e.g., an isotropic Gaussian distribution. Here,  $\mathbf{w}_t$  is the standard Wiener process (or Brownian motion) and  $\mathbf{f}(\mathbf{x}_t, t)$  denotes the drift term that guides the diffusion process, which we will assume to be zero for the most part of this section. Meanwhile,  $g(t)$  represents the diffusion coefficient that controls the noise at each time step  $t \rightarrow T$ . Secondly, the *reverse process* removes the injected noise at each step  $t$  and it is described as

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)]dt + g(t)d\bar{\mathbf{w}}_t, \quad (4.2)$$

where  $\bar{\mathbf{w}}_t$  is the standard Wiener process. To effectively solve Eq. (4.2), one must reliably estimate the score  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$  via training a neural network  $s_\theta(\mathbf{x}, t)$ . The learned weights  $\theta^*$  are obtained using methods for denoising score matching across multiple times steps [42; 43; 45]. The general description of this optimization problem, given by [45], is formulated as

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t, \mathbf{y}, \mathbf{x}_t} [\lambda(t) \|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{y})\|^2], \quad (4.3)$$

where  $t \sim \mathcal{U}(t_0, T)$  is sampled from the uniform distribution  $\mathcal{U}$  over the set of times ranging from a small time  $t_0 \approx 0$  to a larger time  $T$ , while  $\mathbf{y} \sim p(\mathbf{y})$  and  $\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{y})$ . Here,  $p(\mathbf{x}_t|\mathbf{y})$  is the forward process and  $\lambda(t)$  is a positive weighting function.

## 4.2 Diffusion Models from Associative Memories

A fascinating aspect about diffusion models is their process of shaping their energy landscape. Instead of directly learning their energy function  $E_\theta(\mathbf{x}_t, t)$ , diffusion models learn the negative gradient of their energy function or the score function:

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = -\nabla_{\mathbf{x}_t} E_\theta(\mathbf{x}_t, t), \quad (4.4)$$

<sup>1</sup>Assume that  $\mathbf{x}_0 = \mathbf{y} \in \mathbb{R}^D$  are i.i.d samples coming from a data distribution  $p(\mathbf{y})$ .

using the process denoted in Eq. (4.3). However, this particular process does not explain how generalization happens in such models. Instead, it tells a story of diffusion models behaving like AM systems. Specifically, during training, diffusion models are learning how to remove noise from a perturbed memory cue (or query) to obtain a clean memory accordingly to Eq. (4.3). At some point, these models must behave like AM systems where they can effectively recover memories (or stored training data points) from noises. But once a certain threshold (memorization capacity) is exceeded, diffusion models can no longer act like effective denoisers or AM systems, the successive failure in memory recall of these models must facilitate and signal their transition to generative modeling.

Consequently, following the derivation done in [50], we can establish a fundamental connection between diffusion and AM models. Consider the training data distribution in the variance-exploding (VE) setting of  $f(\mathbf{x}_t, t) = 0$  and  $g(t) = \sigma$ . In this case, the marginal probability distribution of new samples can be computed exactly as

$$p(\mathbf{x}_t, t) = \mathbb{E}_{\mathbf{y} \sim \text{data}} \left[ \frac{1}{(2\pi\sigma^2 t)^{\frac{D}{2}}} \exp \left( -\frac{\|\mathbf{x}_t - \mathbf{y}\|_2^2}{2\sigma^2 t} \right) \right]. \quad (4.5)$$

Assuming the empirical distribution of the data  $p(\mathbf{y}) = \frac{1}{K} \sum_{\mu=1}^K \delta^{(D)}(\mathbf{y} - \boldsymbol{\xi}^\mu)$ , where  $\boldsymbol{\xi}^\mu$  represents an individual data point (with data size  $K$ ), this marginal distribution can be written as

$$p(\mathbf{x}_t, t) = \frac{1}{K} \sum_{\mu=1}^K \frac{1}{(2\pi\sigma^2 t)^{\frac{D}{2}}} \exp \left( -\frac{\|\mathbf{x}_t - \boldsymbol{\xi}^\mu\|_2^2}{2\sigma^2 t} \right) \stackrel{\text{def}}{=} \exp \left( -\frac{E^{\text{DM}}(\mathbf{x}_t, t)}{2\sigma^2 t} \right), \quad (4.6)$$

where we also defined the energy  $E^{\text{DM}}$  of diffusion model, which up to state- or  $\mathbf{x}$ - independent terms is equal to

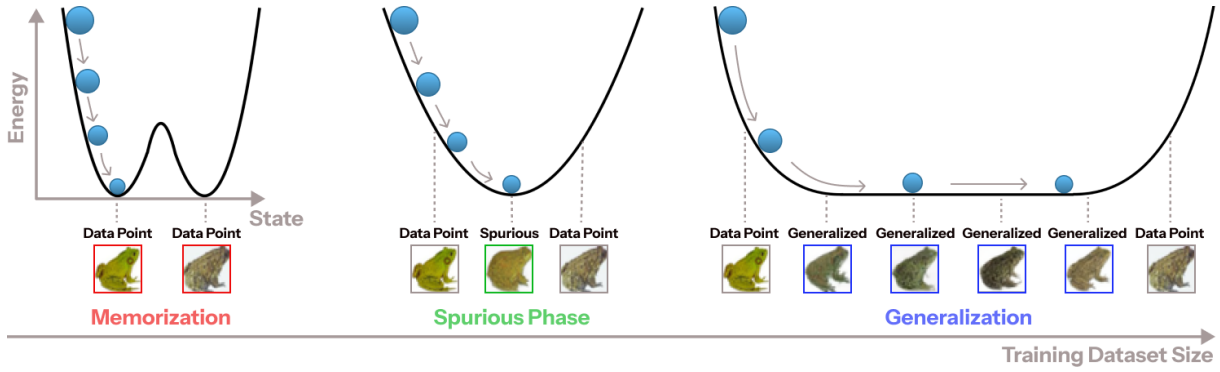
$$E^{\text{DM}}(\mathbf{x}_t, t) = -2\sigma^2 t \log \left[ \sum_{\mu=1}^K \exp \left( -\frac{\|\mathbf{x}_t - \boldsymbol{\xi}^\mu\|_2^2}{2\sigma^2 t} \right) \right]. \quad (4.7)$$

As already observed in [47], the above energy function (4.7) is closely related to that of DenseAMs, which are large memory storage variants of classical Hopfield networks, see Chapter (2).

The core idea behind DenseAMs is to design an energy function that peaks very sharply around the intended memory patterns to prevent the overlapping (or cross talk) between them. Hence, such networks can store and retrieve a much larger number of patterns, compared to the classical Hopfield networks, and scales super-linearly (and possibly exponentially) to the size of the network, allowing the decoupling of information storage capacity from the dimensionality of the data [21; 27]. Of particular interest here is the DenseAM model studied in [60] (see also [31]), which bears strong resemblance to Eq. (4.7):

$$E^{\text{AM}}(\mathbf{x}) = -\beta^{-1} \log \left[ \sum_{\mu=1}^K \exp \left( -\beta \|\mathbf{x} - \boldsymbol{\xi}^\mu\|_2^2 \right) \right], \quad (4.8)$$

where  $\beta$  is the inverse “temperature”, which controls the steepness of the energy landscape around the memories  $\boldsymbol{\xi}^\mu$ .



**Figure 4.2:** A simple illustration depicting the change in the energy landscape as the size of the training dataset is increased. In the small data regime, the diffusion model memorizes the training data points as local minima of the energy. When the amount of training data exceeds the memory capacity of the model, spurious patterns are formed and training data points are no longer energy minima. Subsequent increase of the training set size leads to the generalization phase, which is defined by the formation of continuous manifold of the low energy states. Figure is obtained from [50].

Notice that Eqs. (4.7) and (4.8) are identical if  $\beta = 1/(2\sigma^2 t)$ . The two systems described above have important differences and similarities. In typical AM tasks, the inverse temperature  $\beta$  is kept constant (and is typically large). At the same time, the diffusion energy  $E^{\text{DM}}(\mathbf{x}_t, t)$  describes an intrinsically non-equilibrium system, since the effective temperature explicitly depends on time. However, notice that since the reverse process (4.2) is guaranteed to invert the forward step (4.1), the fixed points of the denoising process are guaranteed to coincide with the original data points. Specifically, both of the above energy functions, Eqs. (4.7) and (4.8), express a competition among the stored data points  $\mathbf{y} \sim p(\mathbf{y})$  to see which one is closer to the query  $\mathbf{x}_t$  at time  $t$ . Hence, although there might be differences in dynamical trajectories for  $E^{\text{DM}}$  and  $E^{\text{AM}}$ , their fixed points must be the same<sup>2</sup>.

Specifically, the manifolds of the data in diffusion models must emerge from the point-like memory storing systems, like AMs, in the limit when they are overloaded with amount of data above the critical memory capacity. In this regime, distinct basins of attraction corresponding to separate memories merge, forming the manifolds of the data. At the boundary of this transition, a separate “phase” corresponding to spurious states, which is ubiquitous in AMs around the *critical memory load*, appears and signals the onset of generalization, see Fig. (4.2) for the simple illustration of this memorization-generalization transition. It is worth noting that DenseAMs typically have an exponentially large memory storage capacity (in the number of neurons  $D$ ) for uncorrelated patterns. However, in the cases of real data, due to the high correlation of samples, the *critical memory load* is much lower than the exponentially large capacity of uncorrelated data — a well-known fact in associative memories [61; 62; 63; 64; 65; 66].

<sup>2</sup>We remind the reader that the fixed points retrieved from the reverse process correspond to  $t = 0$ .

### 4.3 Memorization - Spurious - Generalization Transition

To better illustrate the connection between diffusion models and DenseAMs, we can investigate a simple 2-dimensional toy model, see Fig. (4.3), that exhibits many aspects of the memorization-generalization transition in these two types of models. Specifically, imagine that the training data lies on a unit circle. We are interested in exploring how the shape of the energy function (4.8) changes as the number of training data points, used in training a diffusion model, increases.

Specifically, in the trivial case of a single training point ( $K = 1$ ), there exists only a single memory  $\xi^1$  on the energy landscape of Eq. (4.8), making it independent from the inverse temperature or sharpness value  $\beta$ . In contrast, when there exists two training points ( $K = 2$ ), there exists two corresponding minima (or memories)  $\xi^1$  and  $\xi^2$ :

$$E^{\text{AM}}(\mathbf{x}) = -\beta^{-1} \log \left[ \exp \left( -\beta \|\mathbf{x} - \xi^1\|_2^2 \right) + \exp \left( -\beta \|\mathbf{x} - \xi^2\|_2^2 \right) \right], \quad (4.9)$$

when  $\beta \rightarrow \infty$ . However, for finite values of  $\beta$ , there exists a configuration which yields a minimum:

$$\boldsymbol{\eta} = \arg \min_{\mathbf{x}} E^{\text{AM}}(\mathbf{x}), \quad (4.10)$$

that does not correspond with any of the two training data points or stored patterns, i.e.,  $\boldsymbol{\eta} \neq \xi^1$  and  $\boldsymbol{\eta} \neq \xi^2$ . This “novel” local minimum of the energy is the *spurious state* illustrated in the cartoonish Fig. (4.2).

When the training data size  $K \rightarrow \infty$ , the empirical data distribution of the toy model can be described as a continuous density of states:

$$p(\mathbf{y}) = \frac{1}{\pi} \delta(y_1^2 + y_2^2 - 1). \quad (4.11)$$

The probability of the generated data is proportional (up to terms independent of the state  $\mathbf{x}$ ) to

$$p(\mathbf{x}) \sim \int_{-\infty}^{+\infty} dy_1 dy_2 p(\mathbf{y}) e^{-\beta \|\mathbf{x} - \mathbf{y}\|_2^2} = e^{-\beta(R^2+1)} I_0(2\beta R), \quad (4.12)$$

where  $I_0(\cdot)$  is a modified Bessel function of the first kind<sup>3</sup>. Thus, the energy of the 2D circle model is given by

$$E^{\text{AM}}(R, \phi) = R^2 + 1 - \frac{1}{\beta} \log [I_0(2\beta R)] \underset{\beta \rightarrow \infty}{\approx} (R - 1)^2, \quad (4.13)$$

---

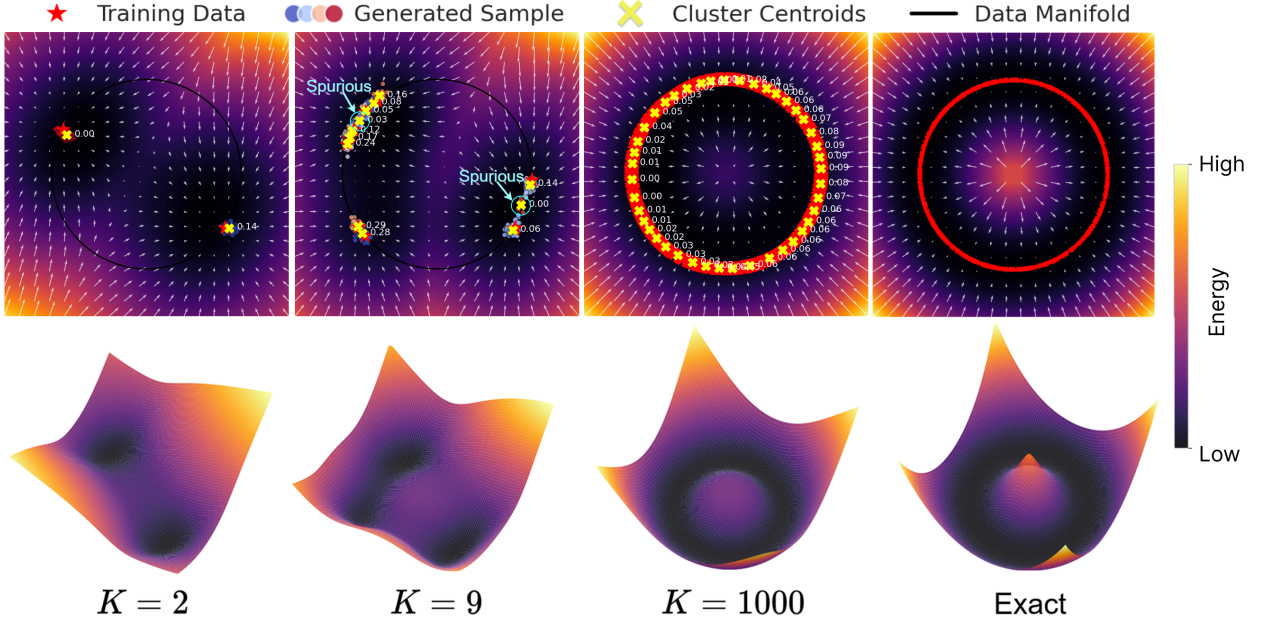
<sup>3</sup>In order to obtain Eq. (4.12), it is easiest to introduce polar coordinates for both the state vector  $\mathbf{x}$  and the training data  $\mathbf{y}$ :

$$\begin{cases} x_1 = R \cos(\phi) \\ x_2 = R \sin(\phi) \end{cases} \quad \begin{cases} y_1 = r \cos(\varphi) \\ y_2 = r \sin(\varphi) \end{cases}$$

The integral (4.12) can then be written as

$$p(\mathbf{x}) \sim \int_0^{2\pi} d\varphi \int_0^\infty r dr \frac{1}{\pi} \delta(r^2 - 1) e^{-\beta[R^2 + r^2 - 2Rr \cos(\varphi - \phi)]} = e^{-\beta(R^2+1)} I_0(2\beta R)$$

and explicitly computed using the definition of the modified Bessel functions [67].



**Figure 4.3:** Energy landscape evolution for the 2D toy model as training data size  $K$  increases. Models trained at  $K \in \{2, 9, 1000\}$  use standard VE-SDE based diffusion pipeline with training data sampled from the unit circle, shown in black for  $K \in \{2, 9\}$ . Generated samples are shown alongside the learned score field  $s_\theta(\mathbf{x}_t, t)$  done via a neural network, aligned with the negative gradient of the energy Eq. (4.7). Hierarchical clustering identifies structure within the generations, with cluster centroid energies visualized by color and numerical value. The rightmost panel shows the exact solution as  $K \rightarrow \infty$  derived in Eq. (4.13). As  $K$  grows, the model initially memorizes individual data points, forming isolated basins. Around  $K = 9$ , *spurious patterns* emerge — distinct low-energy attractors not present in the data — which mark the onset of generalization. At large  $K$ , the model enters a fully generalized regime, where low-energy states lie on a flat, continuous manifold shown in Fig. (4.2). Top-row figure is retrieved from [50].

where  $R$  is the radius of the unit circle and  $\phi$  is the polar angle. Keep in mind, the dependence on  $\phi$  in Eq. (4.13) disappears for the final result.

For our diffusion model, we consider the following forward process which describes the VE-setting:

$$d\mathbf{x}_t = \sigma d\mathbf{w}_t, \quad (4.14)$$

where the drift term  $\mathbf{f}(\mathbf{x}_t, t) = 0$  and  $\mathbf{w}_t$  is Brownian motion. The corresponding reverse process is described as

$$d\mathbf{x}_t = [-\sigma^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)] dt + \sigma^2 d\mathbf{w}_t, \quad (4.15)$$

where the diffusion coefficient  $g(t) = \sigma$  is fixed as 1, matching the radius of the unit circle. Using these SDEs, we trained a set of SDE-based diffusion models, for  $K \in \{2, 9, 1000\}$ , over the time domain of  $t \in [\epsilon, 1]$  where  $\epsilon = 10^{-5}$ , using the objective (4.3). Then, we visualize the energy landscape of each model and record our results in Fig. (4.3).

As expected, the local minima of the resulting  $E^{\text{AM}}$  in Eq. (4.13) form a *continuous manifold*, corresponding to  $R = 1$ . The data samples from the model in Fig. (4.3) occupy the vicinity of that manifold. This behavior describes the fully generalized phase. In the limit of large  $\beta$ , the energy landscape is described by a parabola centered around  $R = 1$ . For our trained diffusion

model at  $K = 1000$ , we see that the exact energy and approximated energy, obtained from the diffusion model, are very much aligned to one another in Fig. (4.3).

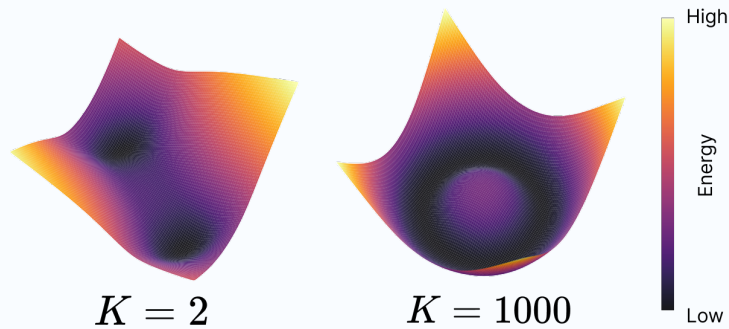
Meanwhile, for small number of data points ( $K = 2$ ), the diffusion model exhibits memorization. The local minima of the energy correspond to the training data points. Importantly, at  $K = 9$ , we are able to observe the first signs of *spurious states*. At this stage, the model begins to learn emergent (different from the training data) local minima of the energy. Subsequent increase of the size of the training set leads to fully generalized behavior, which is illustrated for  $K = 1000$ . At that stage, all of the samples from the model live in close proximity of the exact data manifold. The right panel shows the analytical expression for the energy landscape, defined by Eq. (4.13). Thus, the conventional diffusion modeling pipeline, following Eq. (4.3), agrees very well with the theoretical prediction of the empirical energy (4.13) and the cartoonish illustration in Fig. (4.2).

From the perspective of DenseAMs, we can see a novel phase also exists in diffusion models — *spurious states* — previously overlooked in the memorization-generalization literature of these models [49; 59; 68; 69; 70; 71; 72; 73; 74]. As demonstrated in [50], diffusion models trained on real and high dimensional datasets also follow the same trend illustrated in Fig. (4.3): transitioning from memorization to spurious phase to generalization as the training data size  $K$  increases. Hence, by viewing the problem of generalization as a failure of storing all of the data points as memories, we can provide a novel understanding of the memorization and generalization in generative diffusion models and interestingly, demonstrate the existence of spurious states in such models. This illustrates that diffusion models behave as AM systems in the small data regime, and as generative models in the large data regime.

#### Notebook 4.1: Comparison of diffusion energy and DenseAM energy

In this notebook, we offer the reader the possibility to train a simple diffusion model using data from the 2-D circle as an example. The reader can reconstruct the energy profile of the diffusion model by integrating the score function and compare this energy profile with the energy of DenseAM model.

Checkout the notebook as a [blog post](#), a [colab notebook](#) or as a [raw .ipynb file](#).



## Chapter 5

# Associative Memory: A Machine Learning Model

In this Chapter, we will view Associative Memory networks through the lens of machine learning modeling. After presenting a brief discussion on machine learning modeling, and the sources of error in (machine) learning Section (5.1), we present Associative Memory network as a machine learning model that can be used much like other models in learning, highlighting its inference process, expressivity, application to supervised learning, and its parametric and nonparametric forms, see Section (5.2). Then we discuss how this model can be used for the unsupervised learning task of clustering in Section (5.3). Finally, we elaborate on the connection between Associative Memory and Kernel Machines, and discuss novel Associative Memory models that emerge from this connection in Section (5.4).

### 5.1 Machine Learning Modeling

The purpose of *learning* is to obtain a version of the ground-truth distribution (or equivalently the data-generating function) given (potentially noisy) samples from the ground-truth distribution. The first step is data acquisition. Given data, we choose a model or function class  $\mathcal{F}$  which corresponds to not just a *method* (such as Support Vector Machines [75], Generalized Linear Models [76], Decision trees [77], etc.), but their specific configuration governed by their respective *hyperparameters* such as regularization forms and penalties, trees depth, network architecture and activations, optimization configurations. Given our choice of the function class  $\mathcal{F}$ , the learning process searches (optimizes) for the function  $\hat{f} \in \mathcal{F}$  that (approximately) minimizes the *empirical risk* — the sample loss computed over the training data — or some surrogate of it which better represents the *true risk* — the population loss — or is easier to optimize, such as some continuous version of a discrete loss and/or some form of regularization that mitigates overfitting such as weight penalty or decay, dropout and such.

We currently have an understanding of the factors [78; 79; 80; 81] affecting the *excess risk* of this chosen/learned function — the difference between the true risk of this learned function  $\hat{f}$  and the best possible function  $f^*$ . At a high level, these factor depend on (i) the choice of the function class and its capacity to model the data generating process, (ii) the use of an empirical risk *estimate* instead of the true risk to learn this function, and (iii) the approximation in the empirical risk minimization (ERM) over the class of functions  $g \in \mathcal{F}$ .

For a particular method (decision trees, linear models, neural networks), let  $\mathcal{F}$  denote the function class for some *fixed* hyperparameter  $\lambda \in \Lambda$  (tree depth, number of trees for tree ensembles; regularization parameter for linear and nonlinear models, activation functions, batch size in stochastic gradient descent or SGD, etc.) in the space of valid hyperparameters  $\Lambda$ .

Focusing on supervised learning, for any model or function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  with  $(\mathbf{x}, \mathbf{y}), \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$  generated from a distribution  $p_{\text{data}}$  over  $\mathcal{X} \times \mathcal{Y}$ , and a loss function  $\mathcal{L} : \mathcal{Y} \times \mathcal{Y}$ , the true risk  $R(f)$  and the empirical risk  $R_m(f)$  of  $f$  with  $m$  samples  $\{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^m \sim p_{\text{data}}$  is given by

$$\begin{aligned} R(f) &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}} [\mathcal{L}(\mathbf{y}, f(\mathbf{x}))] = \int \mathcal{L}(\mathbf{y}, f(\mathbf{x})) dp_{\text{data}}, \\ R_m(f) &= \mathbb{E}_m [\mathcal{L}(\mathbf{y}, f(\mathbf{x}))] = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{y}^i, f(\mathbf{x}^i)). \end{aligned} \quad (5.1)$$

We denote the *Bayes optimal model* as  $f^*$  where, for any  $(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}$ ,

$$f^*(\mathbf{x}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbb{E} [\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) | \mathbf{x}]. \quad (5.2)$$

We denote with the following:

$$\bar{f} = \arg \min_{f \in \mathcal{F}} R(f), \quad \hat{f}_m = \arg \min_{f \in \mathcal{F}} R_m(f), \quad (5.3)$$

as the *true risk minimizer*  $\bar{f} \in \mathcal{F}$  and the *empirical risk minimizer*  $\hat{f}_m \in \mathcal{F}$  (with  $m$  samples) in model class  $\mathcal{F}$  respectively.

When performing empirical risk minimization or ERM over  $\mathcal{F}$ , the excess risk is given by

$$\mathcal{E} = R(\hat{f}_m) - R(f^*) = \underbrace{R(\hat{f}_m) - R(\bar{f})}_{\mathcal{E}_{\text{est}}} + \underbrace{R(\bar{f}) - R(f^*)}_{\mathcal{E}_{\text{app}}}, \quad (5.4)$$

which decomposes into two terms: (i) the *approximation risk*  $\mathcal{E}_{\text{app}} = R(\bar{f}) - R(f^*)$ , and (ii) the *estimation risk*  $\mathcal{E}_{\text{est}}(m) = R(\hat{f}_m) - R(\bar{f})$ . For limited number of samples  $m$ , there is a tradeoff between  $\mathcal{E}_{\text{app}}$  and  $\mathcal{E}_{\text{est}}$ , where a larger function class  $\mathcal{F}$  usually reduces  $\mathcal{E}_{\text{app}}$  but increases  $\mathcal{E}_{\text{est}}(m)$  [78; 79]. Roughly speaking, methods are termed *universal approximators* if there is some hyperparameter which ensures that the approximation error  $\mathcal{E}_{\text{app}}$  can be made arbitrarily small. Of course, the flip side is that this can make the corresponding class of functions  $\mathcal{F}$  very large, often increasing the estimation error  $\mathcal{E}_{\text{est}}(m)$  for a fixed  $m$ .

Bottou and Bosquet [80] study the tradeoffs in a “large-scale” setting where the learning is compute bound (in addition to the limited number of samples  $m$ ). Given any computational budget  $T$ , they consider the learning setting “small-scale” when the number of samples  $m$  is small enough to allow for the ERM to be performed to arbitrary precision. In this case, the tradeoff is between the  $\mathcal{E}_{\text{app}}$  and  $\mathcal{E}_{\text{est}}$  terms (as above). They consider the large scale setting where the ERM needs to be approximated given the computational budget and discuss the tradeoffs in the excess risk of an approximate empirical risk minimizer  $\tilde{f}_m \in \mathcal{F}$ . In addition to  $\mathcal{E}_{\text{app}}$  and  $\mathcal{E}_{\text{est}}$ , they introduce the *optimization risk* term  $\mathcal{E}_{\text{opt}} = R(\tilde{f}_m) - R(\hat{f}_m)$  — the excess risk incurred due to

approximate ERM — and argue that, in compute-bound large-scale learning, approximate ERM on all the samples  $m$  can achieve lower excess risk than high precision ERM on a subsample of size  $m' \leq m$ . Fig. (5.1) provides a visual representation of this excess risk decomposition.



**Figure 5.1: Decompositions of excess risk  $\mathcal{E}$ .** We depict the decomposition of  $\mathcal{E}$  incurred by the approximate empirical risk minimizer  $\tilde{f}_m \in \mathcal{F}$  found (usually) with a scalable optimization algorithm in the model class  $\mathcal{F}$  with respect to the Bayes optimal model  $f^*$ . The true risk minimizer  $\bar{f} \in \mathcal{F}$  is the best approximator of the optimal  $f^*$ , while the exact empirical risk minimizer  $\hat{f}_m \in \mathcal{F}$  can be distinct from the true risk minimizer  $\bar{f}$  since we are using the empirical risk (obtained with a finite training set) for our learning instead of the true risk. This figure is partially replicated from Ram et al., 2023 [81]. **Left:** The model class  $\mathcal{F}$  can be such that the Bayes optimal  $f^* \notin \mathcal{F}$ , hence we would have a nonzero approximation risk  $\mathcal{E}_{\text{app}}$  — the difference in the true risk (defined in Eq. (5.1)) between the Bayes optimal  $f^*$  and the true risk minimizer  $\bar{f}$  in our model class  $\mathcal{F}$  — with  $\mathcal{E}_{\text{app}} > 0$ . **Right:** We can also select a large model class  $\mathcal{F}$  such that the Bayes optimal  $f^*$  is in the model class  $\mathcal{F}$  or can be approximated to arbitrary precision with a model in the function class. In this case, the true risk minimizer  $\bar{f} \approx f^*$  will have (almost) zero approximation risk with  $\mathcal{E}_{\text{app}} \approx 0$ . However, it is important to note that the estimation risk  $\mathcal{E}_{\text{est}}$  — the difference in the risk between the true risk minimizer and the empirical risk minimizer — is often related to the size of the model class  $\mathcal{F}$  (for some notion of size), with larger model classes incurring larger estimation risk. The optimization risk  $\mathcal{E}_{\text{opt}}$  — the difference between the risk of the exact and approximate empirical risk minimizers — can also be affected implicitly by the size of the model class  $\mathcal{F}$ , where larger classes require more computational resources for the learning optimization to achieve any specific level of empirical risk approximation; conversely, for fixed computational resources, larger function classes can incur larger optimization risk.

For parametric models, the functions  $f \in \mathcal{F}$  are parameterized with  $\theta \in \Theta$  where  $\mathcal{F} \triangleq \{f_\theta, \theta \in \Theta\}$ , where we explicitly denote the dependence of  $\theta$  in  $f_\theta$ . Thus, the true risk minimizer, and the empirical risk minimizer can be respectively written as:

$$\bar{f} \triangleq f_{\bar{\theta}}, \quad \bar{\theta} = \arg \min_{\theta \in \Theta} R(f_\theta) = \arg \min_{\theta \in \Theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}} [\mathcal{L}(\mathbf{y}, f_\theta(\mathbf{x}))] \quad (5.5)$$

$$\hat{f}_m \triangleq f_{\hat{\theta}_m}, \quad \hat{\theta}_m = \arg \min_{\theta \in \Theta} R_m(f_\theta) = \arg \min_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{y}^i, f_\theta(\mathbf{x}^i)). \quad (5.6)$$

The approximate empirical risk minimizer would be denoted with  $\tilde{f}_{\tilde{\theta}_m}$  with corresponding model parameters  $\tilde{\theta}_m$ . The parameters  $\theta \in \Theta$  corresponds to model/function specific parameters — weights and biases for linear models and neural networks; split dimensions and thresholds, and leaf node values for univariate decision trees. Once these parameters  $\hat{\theta}_m$  are learned from the data, one can make predictions  $f_{\hat{\theta}_m}(\mathbf{x})$  on new inputs  $\mathbf{x} \in \mathcal{X}$  without having to keep the training data  $\{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^m$  around anymore. For nonparametric models such as nearest neighbor models

and (kernel) density estimation based models (such as the Nadaraya-Watson estimator), the training data is also required for making predictions, and thus are considered as part of the “model parameters”.

## 5.2 Associative Memory Network as a Model

The previously discussed Associative Memory networks can be viewed as a parameterized model  $f_{\Xi} : \mathcal{X} \rightarrow \mathcal{X}$  with parameters  $\Xi$ . For the sake of simplicity, let us for now assume that  $\mathcal{X} \subseteq \mathbb{R}^D$ , the  $D$ -dimensional Euclidean space. The interpretation is detailed in the following:

### Model parameters are stored patterns

The model parameters  $\Xi$  can be reshaped as a  $(D \times K)$  matrix, and is usually termed as the  $K$  stored patterns  $\{\xi^\mu \in \mathbb{R}^D, \mu \in [K]\}$  — each stored pattern  $\xi^\mu$  is a  $D$  dimensional vector. Note that these model parameters can be learned.

### Energy function

Given these model parameters  $\Xi$ , we have an energy function of a state  $\mathbf{v} \in \mathcal{X} \subseteq \mathbb{R}^D$ , usually of the following general form:

$$E_\beta(\mathbf{v}; \Xi) = -Q \left( \sum_{\mu=1}^K F(\beta S[\sigma(\mathbf{v}), \xi^\mu]) \right), \quad (5.7)$$

- One can view the  $\mathbf{v}$  as the internal state and  $\sigma(\mathbf{v})$  as its activation. See Section (3.2) of Chapter (3) for a discussion on states and activations.
- The  $S : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  denotes a notion of similarity such as a dot-product or negative squared Euclidean distance.
- $\beta > 0$  is the inverse temperature and controls how much high similarities are magnified and low similarities are diminished. This inverse temperature  $\beta$  controls the sharpness of the energy around the memories, with larger values of  $\beta$  inducing sharper energy landscapes while smaller values generating smoother ones. See Notebook 2.1 in Section (2.3), Chapter (2) for a demonstration.
- The separation function  $F : \mathbb{R} \rightarrow \mathbb{R}$  is a fast-growing function such as  $F(z) = z^p$  or  $F(z) = \exp(z)$  for some  $z \in \mathbb{R}$ .
- The scaling function  $Q : \mathbb{R} \rightarrow \mathbb{R}$  is a monotonic non-decreasing function such as identity  $Q(z) = z$  or logarithm  $Q(z) = \log z$  for some  $z \in \mathbb{R}$ .

### Inference via energy descent

Given a learning rate  $\eta > 0$ , number of steps  $T$  and a clamping mask  $\mathbf{m} \in \{0, 1\}^D$ ,  $f_{\Xi}(\mathbf{x})$  for some  $\mathbf{x} \in \mathcal{X}$  is computed as follows via (clamped) coordinate gradient descent over the

energy function, with the descent initialized as the input  $\mathbf{x}$ :

$$\mathbf{v}^{(0)} \leftarrow \mathbf{x}, \quad (5.8)$$

$$\mathbf{v}^{(t)} \leftarrow \mathbf{v}^{(t-1)} - \eta \mathbf{m} \odot \nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}; \Xi)|_{\mathbf{v}=\mathbf{v}^{(t-1)}}, \quad t \in \llbracket T \rrbracket, \quad (5.9)$$

$$f_{\Xi}^{\mathbf{m}}(\mathbf{x}) \triangleq \mathbf{v}^{(T)}, \quad (5.10)$$

where  $\odot$  denotes the element-wise multiplication of vectors. In the absence of the clamping mask, we drop the  $\mathbf{m}$  subscript and use  $f_{\Xi}$ .

- As the learning rate  $\eta \rightarrow 0$ , this energy gradient descent is defined by the following dynamics:

$$\frac{d\mathbf{v}}{dt} = -\mathbf{m} \odot \nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}; \Xi). \quad (5.11)$$

- The learning rate  $\eta$  does not need to be fixed, and can also vary with time. For example, the learning rate can decay with time.
- When the number of steps  $T$  (also referred to as the number of layers in the Associative Memory network) goes to infinity (that is,  $T \rightarrow \infty$ ), for an appropriately set learning rate  $\eta$  (sufficiently small or appropriately scheduled),  $f_{\Xi}(\mathbf{x})$  will be one of the local minima of the energy function — that is,  $\nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}; \Xi) = \mathbf{0}$ ,  $\nabla_{\mathbf{v}}^2 E_{\beta}(\mathbf{v}; \Xi) \succ 0$  at  $\mathbf{v} = f_{\Xi}(\mathbf{x})$ . These fixed points (local minima) are often termed as the *retrieved memories*. Note that we are seeking local minima and not saddle points which correspond to *meta-stable states* where it is hard to decrease the energy but it is not a local minima. See Demircigil et al., 2017 [27] for a discussion of the energy landscape and the *basins of attraction* for the different local minima.
- The gradient clamping mask  $\mathbf{m}$  enables *clamping* of a subset of the state variables. When  $\mathbf{m}$  is the  $D$ -dimensional all-one vector  $\mathbf{1}_N$ , the complete state vector  $\mathbf{v}$  is modified in Eq. (5.9). If  $\mathbf{m} = [\mathbf{1}_{N'}, \mathbf{0}_{(N-N')}]^{\top}$ , then only the first  $N' < N$  entries of the state vector  $\mathbf{v}$  are allowed to be modified, while the remaining  $(N - N')$  entries of  $\mathbf{v}$  are clamped to their initial values obtained from the input  $\mathbf{x}$ . This clamping and coordinate-wise gradient descent is discussed in Section (3.6) of Chapter (3), see also Fig. (3.4) for a visualization of the clamped energy-descent.
- If the learning rate  $\eta$  is small, and  $T$  is not too large, the input would not be modified significantly, and  $f_{\Xi}(\mathbf{x}) \approx \mathbf{x}$ . If  $\eta$  is large (and not decayed appropriately), we might never arrive at a local minima of the energy.

Given the model parameters  $\Xi$  (that can be interpreted as  $K$  stored patterns), the various hyperparameters — the functions  $Q, F, S$  in the energy function, the inverse temperature  $\beta$ , the learning rate  $\eta > 0$ , the number of layers/steps  $T$  — define the inference process with this model  $f_{\Xi}$ . There is a tight relationship between the energy function and probability density through the [Boltzmann distribution](#) — that is, the density  $p(\mathbf{v})$  of a state  $\mathbf{v}$  is tied its energy  $E(\mathbf{v})$  as  $p(\mathbf{v}) \propto \exp(-E(\mathbf{v}))$ . Given this interpretation, the inference with the described Associative Memory network amounts to a form of likelihood maximization via gradient descent.

**Classic energy for binary patterns**

As an example, if  $\xi^\mu \in \{-1, 1\}^D \forall \mu \in \llbracket K \rrbracket$ ,  $\sigma : \mathbb{R} \rightarrow \{-1, +1\}$ ,  $S[\mathbf{v}, \mathbf{v}'] = \langle v, v' \rangle$ ,  $F(z) = z^2$ ,  $\beta = 1$  and  $Q$  is the identity function, then the corresponding energy function is that of the Classic Hopfield Network (CHN):

$$E(\mathbf{v}; \Xi) = - \sum_{\mu=1}^K (\langle \sigma(\mathbf{v}), \xi^\mu \rangle)^2. \quad (5.12)$$

**Log-sum-exp energy with real valued patterns**

With  $\xi^\mu \in \mathbb{R}^D \forall \mu \in \llbracket K \rrbracket$ ,  $\sigma$  as identity,  $S[\mathbf{v}, \mathbf{v}'] = -1/2 \|\mathbf{v} - \mathbf{v}'\|^2$ ,  $F(z) = \exp(z)$  and  $Q(z) = \log z$ , we obtain the widely used log-sum-exp or LSE energy:

$$E_\beta(\mathbf{v}; \Xi) = - \log \sum_{\mu=1}^K \exp \left( -\beta/2 \|\mathbf{v} - \xi^\mu\|^2 \right). \quad (5.13)$$

Note that common representations of the LSE energy contains a preceding  $(1/\beta)$  term on the right-hand-side to cancel out the  $\beta$ -scaling in the gradient. However, we are removing this here since we only care about the direction of the gradient, and not the magnitude.

**5.2.1 Memory Capacity and Expressivity**

Given the above energy function, we can now consider the set  $\mathcal{M} \subset \mathcal{X}$  of local minima of the energy function in Eq. (5.7) defined as:

$$\mathcal{M} = \{ \mathbf{v} \in \mathcal{X} : \nabla_{\mathbf{v}} E_\beta(\mathbf{v}; \Xi) = 0, \nabla_{\mathbf{v}}^2 E_\beta(\mathbf{v}; \Xi) \succ 0 \}. \quad (5.14)$$

Note that this set of local minima will depend on all but two of the various hyperparameters previously discussed — this set does not depend on the learning rate  $\eta$  and the number of steps  $T$ . In the scenario where the learning rate  $\eta \rightarrow 0$  and the number of layers/steps  $T \rightarrow \infty$ , for any input  $\mathbf{x} \in \mathcal{X}$ , the output  $f_\Xi(\mathbf{x}) \in \mathcal{M}$  as the output is one of the local minima, thereby essentially making  $f_\Xi : \mathcal{X} \rightarrow \mathcal{M}$  a surjective function (many-to-one mapping); note that  $f_\Xi$  for the same  $\Xi$  may not be a surjection with a finite  $T$  especially when the value of  $T$  is small.

For a stored pattern  $\xi^\mu$ , if there exists a local minima  $\mathbf{v} \in \mathcal{M}$  such that  $\mathbf{v} \approx \xi^\mu$ , then the model has *memorized* a stored pattern and is able to approximately retrieve it (given an appropriate input initializing the energy descent). The *memory capacity* of an Associative Memory network is informally defined as the largest number  $K^{\max}$  of randomly generated stored patterns  $\Xi$  such that each stored pattern can be (approximately) retrieved. For example, with the classic energy in Eq. (5.12),  $K^{\max} \sim O(D)$ ; with the log-sum-exp energy in Eq. (5.13), there exists hyperparameters (specifically, values of  $\beta$ ) such that  $K^{\max} \sim O(\exp(D))$ . See Chapter (2) for further discussion of memory capacity.

Given that, under appropriate configurations, the Associative Memory network operates as a surjection  $f_{\Xi} : \mathcal{X} \rightarrow \mathcal{M}$ , the expressivity or the approximation ability of the model  $f_{\Xi}$  is related to the size (cardinality) of  $\mathcal{M}$ . We can increase the cardinality of  $\mathcal{M}$  to up to  $K^{\max}$  by increasing the model size (in terms of the number of model parameters) to up to  $DK^{\max}$  corresponding to  $\Xi$  containing  $K^{\max}$  stored patterns in  $D$ -dimensions. Beyond that point, increasing the number of (randomly generated) stored patterns in  $K^{\max}$  would not increase the cardinality of  $\mathcal{M}$ . There are ways to carefully design the stored patterns such that the cardinality of  $\mathcal{M}$  can go beyond  $K^{\max}$ .

### 5.2.2 Supervised Learning

While we have discussed the Associative Memory network as a model  $f_{\Xi} : \mathcal{X} \rightarrow \mathcal{X}$  that maps from  $\mathcal{X}$  to  $\mathcal{X}$ , in the previously discussed supervised learning setup of Section (5.1), we usually consider models of the form  $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$  mapping from a feature space  $\mathcal{X}$  to an output space  $\mathcal{Y}$ . One way to handle that with an Associative Memory network is to consider a model  $f_{\Xi} : \mathcal{Z} \rightarrow \mathcal{Z}$ , where  $\mathcal{Z} \triangleq \mathcal{X} \times \mathcal{Y}$ . If  $\mathcal{X}$  is a  $d$ -dimensional feature space, and  $\mathcal{Y}$  is a  $k$ -dimensional output space,  $\mathcal{Z}$  would be a  $D = (d + k)$ -dimensional space with the features and output concatenated into the state vector. This concept is also described in Section (3.6) of Chapter (3).

Consider the clamping vector  $\mathbf{m} = [\mathbf{0}_d^{\top} \mathbf{1}_k^{\top}]^{\top} \in \{0, 1\}^D$ , and the matrix  $\mathbf{M} = [\mathbf{0}_{k \times d} \mathbf{I}_k] \in \{0, 1\}^{k \times D}$ . Also consider an uninformative default (potentially learnable) prediction  $\mathbf{y}_0 \in \mathcal{Y}$  — as an example,  $\mathbf{y}_0 = \mathbf{0}_k$  for regression or  $\mathbf{y}_0 = (1/k)\mathbf{1}_k$  for  $k$ -class classification. Then we can define a function  $g_{\Xi} : \mathcal{X} \rightarrow \mathcal{Y}$  mapping features  $\mathbf{x}$  to predictions using parameters  $\Xi$  as follows:

$$\mathbf{v}^{(0)} \leftarrow [\mathbf{x}^{\top}, \mathbf{y}_0^{\top}]^{\top}, \quad (5.15)$$

$$\mathbf{v}^{(t)} \leftarrow \mathbf{v}^{(t-1)} - \eta \mathbf{m} \odot \nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}; \Xi)|_{\mathbf{v}=\mathbf{v}^{(t-1)}}, \quad t \in \llbracket T \rrbracket, \quad (5.16)$$

$$g_{\Xi}(\mathbf{x}) \triangleq \mathbf{M} \mathbf{v}^{(T)}. \quad (5.17)$$

See Fig. (3.4) for a visualization this energy minimization based inference process for a supervised learning problem. There are a few important things to note here:

- The energy  $E_{\beta}(\mathbf{v}; \Xi)$  now depends on a similarity function  $S : \mathcal{Z} \times \mathcal{Z}$ , which can incorporate similarity between the features in  $\mathcal{X}$  and the outputs in  $\mathcal{Y}$ . For example, the similarity function  $S$  can be defined as  $S[\mathbf{z}, \mathbf{z}'] = \lambda S_{\mathcal{X}}[\mathbf{x}, \mathbf{x}'] + (1 - \lambda) S_{\mathcal{Y}}[\mathbf{y}, \mathbf{y}']$  where  $\mathbf{z} = [\mathbf{x}^{\top} \mathbf{y}^{\top}]^{\top}$  ( $\mathbf{z}'$  defined accordingly with  $\mathbf{x}', \mathbf{y}'$ ), and  $S_{\mathcal{X}} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ ,  $S_{\mathcal{Y}} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  are feature and output specific similarity functions.
- During the energy descent, the features values  $\mathbf{x} \in \mathcal{X}$  provide the initialization  $\mathbf{v}^{(0)} = [\mathbf{x}^{\top} \mathbf{y}_0^{\top}]^{\top}$ , and then are not modified at all because of the clamping mask  $\mathbf{m}$ .
- The final output of  $g_{\Xi}$  is obtained by extracting the last  $k$ -dimensions of the final state  $\mathbf{v}^{(T)}$  for the energy descent with the element selecting matrix  $\mathbf{M}$ .
- We are considering a gradient descent over the energy  $E_{\beta}(\mathbf{v}; \Xi)$ , which would define a density over  $(\mathcal{X} \times \mathcal{Y})$ . However, we are clamping the state variables corresponding  $\mathcal{X}$  to the input  $\mathbf{x}$ . This roughly corresponds to a conditional density over  $\mathcal{Y}$ , and the clamped energy descent corresponds to a conditional likelihood maximization.

### 5.2.3 Nonparametric vs Parametric Models

#### Nonparametric Models

An important question with an Associative Memory model  $f_{\Xi}$  is the process of obtaining the model parameters  $\Xi$  (also known as the “stored patterns”). Given a set of patterns  $\{\xi^\mu, \mu \in \llbracket K \rrbracket\}$  (possibly from an unknown distribution  $p_{\text{data}}$ , that is,  $\xi^\mu \sim p_{\text{data}}$ ), we can consider a nonparametric form of the model where  $\Xi$  is just all of the data  $\{\xi^\mu, \mu \in \llbracket K \rrbracket\}$ , with the size of the model (which would be  $O(KD)$  when each stored pattern  $\xi^\mu$  is of size  $D$ ) growing with the number of stored patterns (which is  $K$ ). As discussed previously in Section (5.2.1), the corresponding energy function  $E_\beta(\cdot; \Xi)$  can have up to  $O(\min\{K, K^{\max}\})$  local minima, where  $K^{\max}$  is the capacity of the model. Note that if the stored patterns have a specific structure, then the number of local minima can be higher than  $K^{\max}$ . For the supervised learning setup discussed in Section (5.2.2), the stored patterns  $\xi^\mu$  would be feature-output pairs  $(\mathbf{x}, \mathbf{y})$ ,  $\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$  in the training set.

The single-step retrieval dynamics (that is, number of layers  $T = 1$ ) of the nonparametric Associative Memory networks has recently also been recently interpreted as the solution of a specific nonparametric support vector regression problem [82]. Different choices of kernel functions and training data preprocessing result in different energy functions  $E_\beta(\cdot; \Xi)$ .

#### Parametric Models

One can also consider a parametric form of the model  $f_{\Xi}$ , where the size of  $\Xi$  is pre-specified, and the parameters are learned using the data. As an example, we can say that the size of  $\Xi$  is such that it can store only  $K$  patterns  $\xi^\mu, \mu \in \llbracket K \rrbracket$  of size  $D$ . However, here we are allowed to learn these patterns  $\xi^\mu$ . Given a set of training examples  $S = \{\mathbf{z}^i\}_{i=1}^m$ , a loss function  $\mathcal{L}$  and a regularizer  $R$ , at a high level, we can learn  $\Xi$  by solving the following problem:

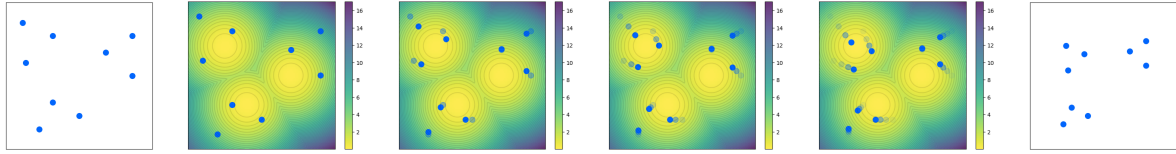
$$\min_{\Xi} R(\Xi) + \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{z}^i, f_{\Xi}(\mathbf{z}^i)). \quad (5.18)$$

If the training example  $\mathbf{z}_i$  is a feature-label pair  $(\mathbf{x}_i, \mathbf{y}_i)$ , and the loss is the negative cross-entropy loss between the labels, and the Associative Memory model is as defined in Section (5.2.2), then  $\mathcal{L}(\mathbf{z}^i, f_{\Xi}(\mathbf{z}^i))$  would simplify to  $\text{NegativeCrossEntropy}(\mathbf{y}^i, f_{\Xi}(\mathbf{x}^i))$  as in a standard classification problem. The regularization  $R(\Xi)$  can be utilized to avoid overfitting. For example,  $R(\Xi) = \lambda \|\Xi\|_{2,1} = \lambda \sum_{\mu \in \llbracket K \rrbracket} \|\xi^\mu\|$  penalizes the norm of the learnable stored patterns, scaled with a  $\lambda \in \mathbb{R}_+$ . One can also consider a regularization that enforces the learnable stored patterns to be well separated so as to make the memory retrieval process more efficient and robust [83; 84]. The loss  $\mathcal{L}$  and the regularization  $R$  can be specified in a problem dependent manner. As the model  $f_{\Xi}$  corresponds to a  $T$ -layer recursive network, we can learn  $\Xi$  with (stochastic) gradient descent given that the loss  $\mathcal{L}$  and the regularization  $R$  are differentiable.

## 5.3 Clustering

Consider an Associative Memory network based model with parameters  $\Xi = \{\xi^\mu \in \mathbb{R}^D, \mu \in \llbracket K \rrbracket\}$  and an energy function  $E_\beta(\cdot; \Xi) : \mathbb{R}^D \rightarrow \mathbb{R}$  as defined in Section (5.2). With any input  $\mathbf{x} \in \mathbb{R}^D$ ,

the energy descent involved in the model output  $f_{\Xi}(\mathbf{x})$  intuitively would move  $\mathbf{x}$  towards the local energy minimum closest to it. When the number of local minima is relatively small, for all input, the outputs will contract towards this small set of local minima. An example of this behaviour is shown in Fig. (5.2). One way of thinking about this contraction effect is the following — the Associative Memory model moves relatively close-by points closer together, potentially moving far-away points even farther. What is interesting is that this is a collective contraction effect on the whole set of inputs even though the model operates on all the input independently.



**Figure 5.2: Contraction of points over the energy landscape.** Here we demonstrate how the Associative Memory model can contract a set of points via gradient descent over the energy landscape. **Column 1:** The initial set of 9 points. **Column 2:** The initial set of points are overlaid on an energy landscape with 3 local minima — lighter colour denotes lower energy. **Columns 3-5:** Each point separately undergoes the energy descent for 3 steps corresponding to a 3-layer model. **Column 6:** The outputs of the model applied separately to each point in the set are now a contracted version of the initial set (column 1).

If the set of input points are all close-by to begin with, and the energy local minima are relatively more spread out, then all points could potentially contract towards the same local minima. However, if the input points are as spread out as the energy local minima, then the contraction effect would lead to input points getting more *clustered* — a subset of the points getting closer to each other while each subset getting farther away from each other. This capability of the Associative Memory network makes it quite useful for the classical problem of clustering.

### 5.3.1 Euclidean Clustering

Given a set of points  $S = \{\mathbf{x}^i \in \mathbb{R}^D, i \in \llbracket m \rrbracket\}$  in a  $D$ -dimensional Euclidean space, a commonly studied clustering problem is the  $k$ -means clustering problem, which seeks to solve the following discrete optimization problem:

$$\min_{\mathbf{c}^1, \dots, \mathbf{c}^k \in \mathbb{R}^D} \sum_{i=1}^m \min_{j \in \llbracket k \rrbracket} \|\mathbf{x}^i - \mathbf{c}^j\|^2. \quad (5.19)$$

This problem seeks to partition the set of points  $S$  into  $k$  disjoint subsets  $C^j, j \in \llbracket k \rrbracket$ , with a prototype or center  $\mathbf{c}^j \in \mathbb{R}^d$  for each subset  $C^j$ , ensuring that squared Euclidean distance between a point in the subset and the corresponding center is small. This is a NP-hard problem even for  $k = 2$  [85], and Lloyd’s algorithm [86] is the most commonly used approximate algorithm though many more efficient algorithms with improved approximation guarantees have been developed. The hardness of this problem is partially due to the discrete nature of the objective in Eq. (5.19), and thus usually requires discrete algorithms. This objective in its current form is not conducive to gradient descent based solutions prevalent in modern machine learning. One can modify this discrete objective into a continuous one by replacing the  $\min_{j \in \llbracket k \rrbracket}$  in Eq. (5.19) with a soft-min

function, leading to soft or fuzzy  $k$ -means clustering [87; 88]:

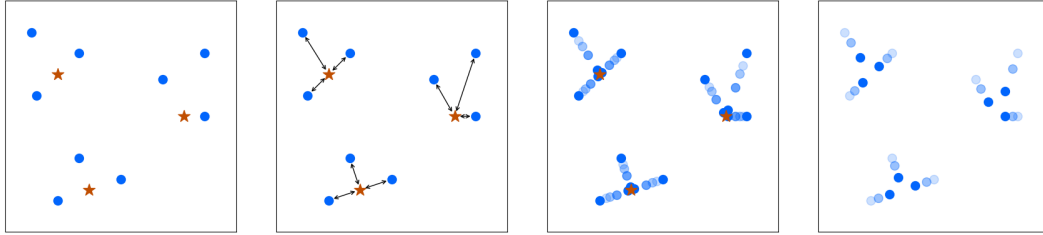
$$\min_{\mathbf{c}^1, \dots, \mathbf{c}^k \in \mathbb{R}^D} \sum_{i=1}^m \sum_{j \in [k]} \frac{\exp(-\gamma \|\mathbf{x}^i - \mathbf{c}^j\|^2) \|\mathbf{x}^i - \mathbf{c}^j\|^2}{\sum_{j' \in [k]} \exp(-\gamma \|\mathbf{x}^i - \mathbf{c}^{j'}\|^2)}, \quad (5.20)$$

where  $\gamma > 0$  is a hyperparameter. The above objective is an upper bound of the  $k$ -means objective in Eq. (5.19) and we would be minimizing the upper bound. Larger values of  $\gamma$  make the upper bound tighter.

Instead of relaxing the discrete assignments of points to clusters  $\min_{j \in [k]} \|\mathbf{x}^i - \mathbf{c}^j\|^2$ , one can emulate the discrete assignments by “moving” a point  $\mathbf{x}^i$  to its closest cluster  $\mathbf{c}^{j^*(\mathbf{x}^i)}$  with  $j^*(\mathbf{x}^i) \triangleq \arg \min_{j \in [k]} \|\mathbf{x}^i - \mathbf{c}^j\|^2$  using the contraction capability of Associative Memory networks and using the term  $\|\mathbf{x}^i - \mathbf{c}^{j^*(\mathbf{x}^i)}\|^2$ , the amount by which the point was “moved” [60]. Thus, the  $k$ -means objective in Eq. (5.19) can be re-written as:

$$\begin{aligned} \min_{\mathbf{c}^1, \dots, \mathbf{c}^k \in \mathbb{R}^D} \sum_{i=1}^m \min_{j \in [k]} \|\mathbf{x}^i - \mathbf{c}^j\|^2 &\equiv \min_{\mathbf{c}^1, \dots, \mathbf{c}^k \in \mathbb{R}^D} \sum_{i=1}^m \|\mathbf{x}^i - \mathbf{c}^{j^*(\mathbf{x}^i)}\|^2 \\ &\equiv \min_{\Xi} \sum_{i=1}^m \|\mathbf{x}^i - f_{\Xi}(\mathbf{x}^i)\|^2 \text{ if } f_{\Xi}(\mathbf{x}^i) \approx \mathbf{c}^{j^*(\mathbf{x}^i)}. \end{aligned} \quad (5.21)$$

This distinction and equivalence is visualized in Fig. (5.3).



**Figure 5.3: Computing the  $k$ -means objective with points and cluster centers.** The computation of the  $k$ -means objective requires us to implicitly or explicitly assign points to clusters. **Column 1:** We are given a set of points  $\bullet$  and centers  $\star$ . **Column 2:** The  $k$ -means objective in Eq. (5.19) assigns each point to its closest center ( $\bullet \leftrightarrow \star$ ), and then sums this distance-to-closest-center over all points. **Column 3:** Instead of relaxing the discreteness in the  $k$ -means objective as in soft  $k$ -means in Eq. (5.20), CIAM [60] uses an Associative Memory network with (learnable) parameters  $\Xi$  to effectively relocate each point to its closest center, and then considers the sum of these per-point-relocation in Eq. (5.21) as a surrogate for the  $k$ -means objective. **Column 4:** Instead of complete contraction to the cluster centers, sometimes it might be beneficial to partially contract to the cluster centers [89].

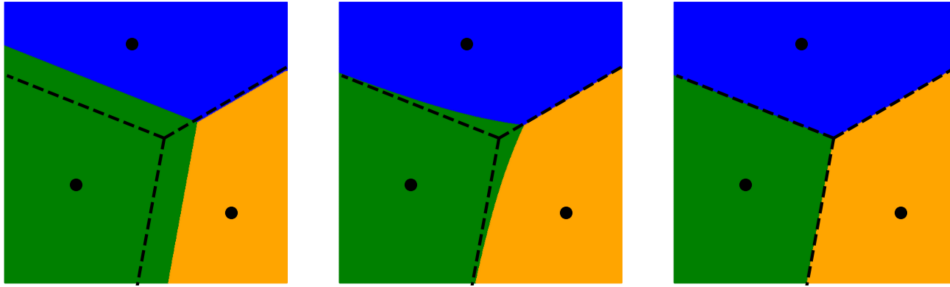
What we need then is an Associative Memory model  $f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$  such that  $f_{\Xi}(\mathbf{x}^i) \approx \mathbf{c}^{j^*(\mathbf{x}^i)}$ . If use the  $k$  cluster centers as the stored patterns, that is  $\Xi \triangleq \{\mathbf{c}^1, \dots, \mathbf{c}^k\}$ , then the condition  $f_{\Xi}(\mathbf{x}^i) \approx \mathbf{c}^{j^*(\mathbf{x}^i)}$  would be satisfied if the basins of attraction around each cluster center (which is also the stored pattern) matches the *Voronoi partition* of the input space given the cluster centers. Given  $k$  centers, the Voronoi partition of the space is (i) a  $k$ -partition of space with each partition corresponding to a specific center, and (ii) any point in a specific partition has its corresponding center as its closest center. Given 3 centers, Fig. (5.4) shows the Voronoi partition

of the input space.

With  $\Xi \triangleq \{\mathbf{c}^1, \dots, \mathbf{c}^k\}$ , and the following energy function  $E_\beta(\cdot; \Xi)$  for an appropriately large inverse temperature  $\beta > 0$  and number of layers  $T$ , we can get the desired behaviour:

$$E_\beta(\mathbf{v}; \Xi) = -\frac{1}{\beta} \log \sum_{j=1}^k \exp(-\beta \|\mathbf{v} - \mathbf{c}^j\|^2). \quad (5.22)$$

The desired behavior of having  $f_\Xi(\mathbf{x}^i) \approx \mathbf{c}^{j^*(\mathbf{x}^i)}$  corresponds to the basins of attraction of each stored pattern  $\xi^\mu, \mu \in \llbracket k \rrbracket$  matching Voronoi partition of the space given the memories/centers  $\{\xi^\mu = \mathbf{c}^\mu, \mu \in \llbracket k \rrbracket\}$ . The dependence on  $\beta$  for  $T = 10$  layers of the DenseAM is visualized in Fig. (5.4). This allows us to solve the discrete clustering problem in Eq. (5.19) with the

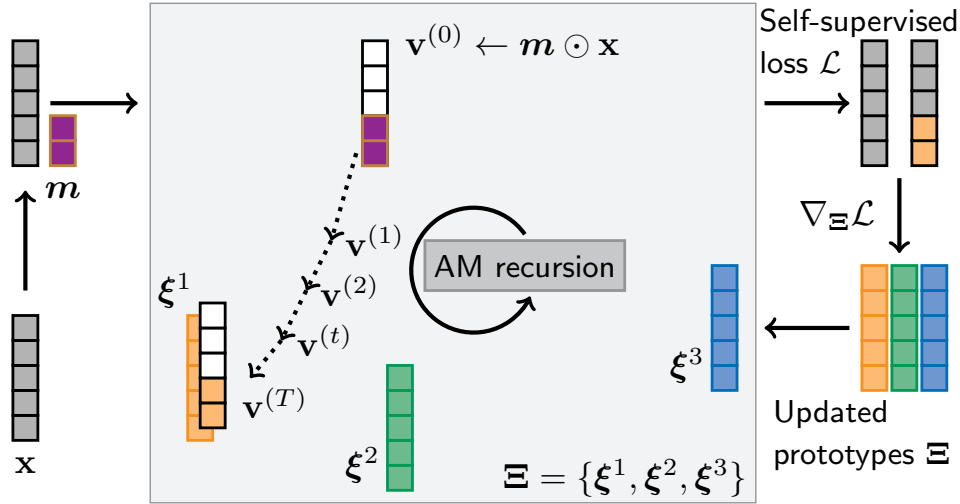


**Figure 5.4: Basins of attraction vs Voronoi partition.** The basins of attraction of the given memories/centers (black dots  $\bullet$ ) for different  $\beta$  values with a 10-layer Associative Memory network  $f_\Xi$  ( $T = 10$ ) are shown by the colored regions. Dashed lines show the desired Voronoi partition. As the value of  $\beta$  increases, the basins of attraction start aligning with the desired Voronoi partitions. **Column 1:** For a small inverse-temperature  $\beta = 0.001$ , the Voronoi partition does not align with the basins of attraction of the Associative Memory network. **Column 2:** With a higher inverse-temperature  $\beta = 10$ , the basins of attraction partially align with the Voronoi partition. **Column 3:** With a high inverse-temperature  $\beta = 100$ , the basins of attraction and the Voronoi partition are practically indistinguishable. This figure is replicated from Saha et al., 2023 [60].

re-written objective (5.21) completely with gradient descent, since we can differentiate through the  $T$  layers of the Associative Memory network. Additionally, we can leverage the clamped inference procedure in Associative Memory networks to extend the standard clustering objective over the training set  $S$  to effectively give us a *self-supervised clustering loss* by creating multiple versions of each point  $\mathbf{x} \in S$  — thereby enabling self-supervision in clustering. Concretely, we can mask each input  $\mathbf{x}$  with a random  $\mathbf{m} \in \{0, 1\}^d$  to form the input  $\mathbf{m} \odot \mathbf{x}$  to the Associative Memory network. Then, we can perform clamped inference  $f_\Xi^{\bar{\mathbf{m}}}(\mathbf{m} \odot \mathbf{x})$  with the clamping mask  $\bar{\mathbf{m}}$  which is the complement of the mask  $\mathbf{m}$ . Thus, our standard clustering loss in Eq. (5.21) would be extended as following from the standard clustering loss on the left to the self-supervised clustering loss on the right:

$$\min_{\Xi} \sum_{\mathbf{x} \in S} \|\mathbf{x} - f_\Xi(\mathbf{x})\|^2 \xrightarrow{\text{self-supervision}} \min_{\Xi} \sum_{\mathbf{x} \in S} \mathbb{E}_{\mathbf{m}} \|\mathbf{m} \odot (\mathbf{x} - f_\Xi^{\bar{\mathbf{m}}}(\mathbf{m} \odot \mathbf{x}))\|^2 \quad (5.23)$$

The overall clustering process, which involves the learning of the stored patterns  $\Xi$ , is visualized in Fig. (5.5).



**Figure 5.5: Euclidean clustering with DenseAM.** For  $\mathbf{x} \in S$ , we first apply a mask (in purple)  $\mathbf{m} \in \{0, 1\}^D$  to  $\mathbf{x}$  to get the initial iterate  $\mathbf{v}^{(0)}$  for the AM recursion. With  $T$  recursions, we have a completed version  $\mathbf{v}^{(T)}$ . The use of the mask  $\mathbf{m}$  is optional, and allows for a semi-supervised clustering loss by leveraging the clamped inference  $f_{\Xi}^{\mathbf{m}}(\mathbf{x})$  in Associative Memory networks; see Eq. (5.23) and Saha et al., 2023 [60, Section 3.4] for details. In the limiting case, there is no mask (that is,  $\mathbf{m} = \mathbf{1}_D$ ) and  $\mathbf{v}^{(0)} \leftarrow \mathbf{x}$  and we do the unclamped inference  $f_{\Xi}(\mathbf{x})$ . The stored patterns  $\Xi$  are updated with the gradient  $\nabla_{\Xi} \mathcal{L}$  on the loss in Eq. (5.21). This figure is replicated from Saha et al., 2023 [60].

### 5.3.2 Deep Clustering

For data modalities, such as image or text, it is often necessary to first learn information-preserving Euclidean representations before clustering these learned representations. This problem of jointly learning representations and clustering is often referred to as *deep clustering* [90; 91; 92]. One common way to learn information-preserving representations is to use an auto-encoder and minimize the reconstruction error, where  $e_{\phi} : \mathcal{X} \rightarrow \mathbb{R}^D$  is a domain-specific encoder (parameterized with  $\phi$ ) that maps the input (images, text) into a latent Euclidean space, which is then used to reconstruct the original data using a decoder  $d_{\vartheta} : \mathbb{R}^D \rightarrow \mathcal{X}$  (parameterized with  $\vartheta$ ) which often mirrors the domain-specific encoder, and the reconstruction loss defined as  $\mathcal{L}_r(\mathbf{x}, d_{\vartheta}(e_{\phi}(\mathbf{x})))$  for a loss  $\mathcal{L}_r : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , giving us the following learning problem given a dataset  $S \subset \mathcal{X}$ :

$$\min_{\phi, \vartheta} \sum_{\mathbf{x} \in S} \mathcal{L}_r(\mathbf{x}, d_{\vartheta}(e_{\phi}(\mathbf{x}))) \quad (5.24)$$

A simple but useful baseline is to just solve the above problem, and then perform  $k$ -means clustering on latent representations  $e_{\phi}(S) = \{e_{\phi}(\mathbf{x}), \mathbf{x} \in S\}$ . However, as we are already learning representations, it is beneficial to steer the learned representations to already have a favourable clustered structure.

This is often obtained by augmenting the reconstruction loss  $\mathcal{L}_r(\mathbf{x}, d_{\vartheta}(e_{\phi}(\mathbf{x})))$  with some form of a clustering loss  $\mathcal{L}_c(e_{\phi}(\mathbf{x}), \{\mathbf{c}^1, \dots, \mathbf{c}^k\})$ , like a relaxed version of  $\min_{j \in [k]} \|e_{\phi}(\mathbf{x}) - \mathbf{c}^j\|^2$ , where  $\mathbf{c}^j \in \mathbb{R}^D, j \in [K]$  are learnable cluster centers in the latent space; see Eq. (5.20) as an example of a continuous clustering loss. Thus, the overall learning problem can be written as following for a

regularization hyperparameter  $\lambda \in [0, 1]$ :

$$\min_{\phi, \vartheta, \{\mathbf{c}^1, \dots, \mathbf{c}^k\} \subset \mathbb{R}^D} \sum_{\mathbf{x} \in S} (1 - \lambda) \mathcal{L}_r(\mathbf{x}, d_{\vartheta}(e_{\phi}(\mathbf{x}))) + \lambda \mathcal{L}_c(e_{\phi}(\mathbf{x}), \{\mathbf{c}^1, \dots, \mathbf{c}^k\}). \quad (5.25)$$

The hyperparameter  $\lambda$  balances the reconstruction and clustering loss as there is an inherent tradeoff between (i) preserving information in the latent space with  $e_{\phi}(\mathbf{x}) \not\approx e_{\phi}(\mathbf{x}')$  for  $\mathbf{x} \neq \mathbf{x}'$  — thus low reconstruction loss, and (ii) forming tight clusters in the latent space by mapping all points within a cluster to almost the same representation, that is  $e_{\phi}(\mathbf{x}) \approx \mathbf{c}^{j^*(\mathbf{x})}$  where  $j^*(\mathbf{x}) \triangleq \arg \min_{j \in [k]} \|e_{\phi}(\mathbf{x}) - \mathbf{c}^j\|^2$  — giving us low clustering loss. The goal is to find the sweet spot, which allows us to have low reconstruction loss (preserving necessary information), while forming a clustered structure in the latent space by pushing both the representations and the cluster centers to have low clustering loss. There is also an implicit objective of forming well-balanced clusters and avoiding *representation collapse*, where all points end up in the same cluster in the latent space.

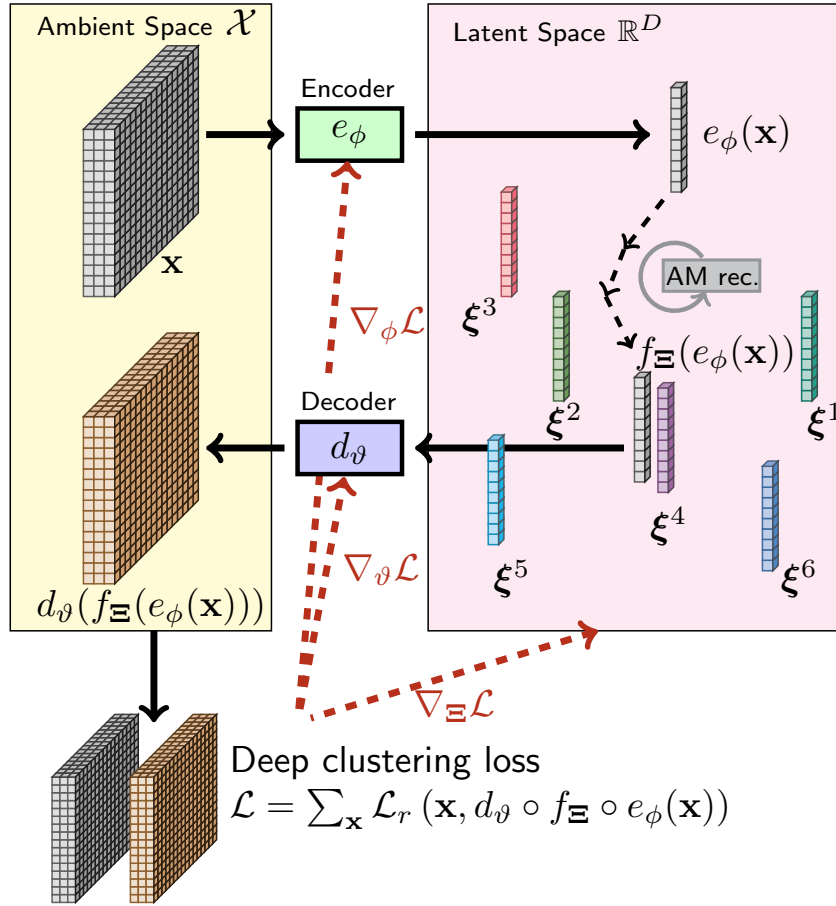
By viewing Associative Memory network as a contractive layer (see Fig. (5.2)), we can introduce a clustered structure in the learnable latent space in an alternate manner. Instead of maintaining an objective  $\mathcal{L}_c$  for clustering that pushes the latent representations to be clustered when minimized, we can employ the contractive nature of the Associative Memory networks to directly have a clustered structure in the latent space, and just focus on optimizing the reconstruction loss of this structured latent space by contracting the latent space before reconstructing [89].

Given a domain-specific encoder  $e_{\phi}$ , a corresponding decoder  $d_{\vartheta}$ , and an Associative Memory model  $f_{\Xi}$  serving as a contraction layer, we can solve the following optimization problem:

$$\min_{\phi, \vartheta, \Xi} \sum_{\mathbf{x} \in S} \mathcal{L}_r(\mathbf{x}, d_{\vartheta}(f_{\Xi}(e_{\phi}(\mathbf{x}))). \quad (5.26)$$

Here, the input  $\mathbf{x}$  is first encoded to the latent space as  $e_{\phi}(\mathbf{x}) \in \mathbb{R}^D$ , and passed through the contraction layer to get  $f_{\Xi}(e_{\phi}(\mathbf{x}))$ . Then the original input is reconstructed with the decoder to get  $d_{\vartheta}(f_{\Xi}(e_{\phi}(\mathbf{x})))$ . In contrast to the use of Associative Memory networks in vanilla Euclidean clustering [60] where we want complete contraction — the points are relocated to the closest cluster center — in this setup, it is beneficial to only consider partial contraction — the points are modified to have a more clustered structure, but the output of the model  $f_{\Xi}$  are still distinct for distinct models. This is visualized in Fig. (5.3) with Fig. (5.3) (Column 3) showing complete contraction, while Fig. (5.3) (Column 4) visualizing partial contraction.

The overall learning procedure is shown in Fig. (5.6). This method provides a single loss function that simultaneously ensures that the loss of information is minimized, while pushing the latent representations to have a clustered structure through the Associative Memory network; no separate clustering loss is utilized here. This loss function, and thus the resulting deep clustering scheme, is agnostic to the data modality (images or text or something else) and the corresponding encoder and decoder architectures. Given that the Associative Memory network model  $f_{\Xi}$  is differentiable, with respect to its learnable parameters  $\Xi$  and its output, we can perform deep clustering by the solving Eq. (5.26) with (stochastic) gradient descent provided the encoder  $e_{\phi}$  and decoder  $d_{\vartheta}$  are



**Figure 5.6: Deep clustering with DenseAM.** Given an input  $\mathbf{x} \in \mathcal{X}$  in the ambient space  $\mathcal{X}$ , the encoder  $e_\phi$  maps  $\mathbf{x}$  to the latent space to get  $e_\phi(\mathbf{x}) \in \mathbb{R}^D$ . Then we use the (partial) contraction capability of a Associative Memory model  $f_\Xi$  to move the latent representation from  $e_\phi(\mathbf{x})$  to  $f_\Xi \circ e_\phi(\mathbf{x})$  towards one of the memories. This contracted representation is then mapped back to the ambient space  $\mathcal{X}$  with the decoder  $d_\vartheta$  to get  $d_\vartheta \circ f_\Xi \circ e_\phi(\mathbf{x})$ . For the purposes of learning the encoder, decoder and Associative Memory network parameters, we utilize the reconstruction loss  $\mathcal{L}_r(\mathbf{x}, d_\vartheta \circ f_\Xi \circ e_\phi(\mathbf{x}))$  and backpropagate the gradients with respect to the parameters  $\phi, \Xi, \vartheta$ . The solid arrows denote the forward-pass  $\mathbf{x} \rightarrow e_\phi(\mathbf{x}) \rightarrow f_\Xi(\mathbf{x}) \rightarrow d_\vartheta \circ f_\Xi \circ e_\phi(\mathbf{x})$  to compute the single loss term in Eq. (5.26). The dashed arrows denote the backward pass showing the single loss driving all updates. This figure is replicated from Saha et al., 2024 [89].

differentiable with respect to their respective parameters and output.

## 5.4 Kernel Machines

Revisiting the energy function of an Associative Memory network  $f_\Xi$  with parameters  $\Xi$  in Eq. (5.7), and denoting the  $F(\beta S[\mathbf{v}, \xi^\mu])$  term with  $\kappa(\mathbf{v}, \xi^\mu)$ , we can write the energy function as:

$$E_\beta(\mathbf{v}; \Xi) = -Q \left( \sum_{\mu \in [K]} \kappa(\mathbf{v}, \xi^\mu) \right), \quad (5.27)$$

where the  $\sum_\mu \kappa(\mathbf{v}, \xi^\mu)$  term can be interpreted as a *kernel sum* with the kernel  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , the core computation in kernel machines [93]. This simple observation allows us to leverage the

rich literature on kernel machine for the development of novel Associative Memory networks with unique capabilities. Two main areas of research in kernel machines focus on the following:

- A lot of research focused on the development and use of expressive domain-specific kernels  $\kappa$ , and the understanding of their properties such as expressivity and generalization. In the context of Associative Memory networks, this corresponds to the development of novel domain-specific energy functions, since one can create an energy given a kernel function through Eq. (5.27), thereby expanding the applicability of these models to new domains.
- Every inference in vanilla kernel machines requires the computation of the kernel sum  $\sum_{\mu} \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu})$  which (i) implies that we need to keep the set  $\{\boldsymbol{\xi}^{\mu}\}_{\mu \in [K]}$  even for inference, and (ii) leads to extremely expensive training and inference as each inference is naively  $O(K)$ , the number of memories (or terms in the kernel sum). A lot of research focused on improving the computational time and space complexity of these kernel-sum computations. This corresponds to improving the computational time and space complexity of the computation of the energy and thus the energy gradient — this would speed up each energy descent step and thus the overall inference with a model  $f_{\Xi}$ .

#### 5.4.1 Random Features

Roughly speaking, for a [symmetric positive definite kernel](#)  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , there exists an implicit feature map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$ , where  $\mathcal{H}$  is a Reproducing Kernel Hilbert space, such that for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ ,  $\kappa(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ . If an explicit feature map  $\phi$  is available, a kernel sum can be simplified as follows:

$$\sum_{\mu=1}^K \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu}) = \sum_{\mu=1}^K \langle \phi(\mathbf{v}), \phi(\boldsymbol{\xi}^{\mu}) \rangle = \left\langle \phi(\mathbf{v}), \sum_{\mu=1}^K \phi(\boldsymbol{\xi}^{\mu}) \right\rangle, \quad (5.28)$$

where we would just need to compute the  $\sum_{\mu=1}^K \phi(\boldsymbol{\xi}^{\mu})$  term once, and use it for any subsequent inference, thereby removing the  $O(K)$  dependence both from the time and space complexity of the inference. However, note that the computational complexities now depends on the size of the feature map  $\phi(\mathbf{v})$  and  $\phi(\boldsymbol{\xi}^{\mu})$ .

A commonly used and expressive kernel is the [RBF \(radial basis function\) kernel](#)  $\kappa : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}_+$  with  $\kappa(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$  for a scaling parameter  $\gamma > 0$ . This kernel possesses an infinite dimensional feature map  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{\infty}$ , thus making the explicit feature map practically unusable. Various indexing schemes have been developed and analysed [94; 95; 96] to speed up the computation of kernel sums..

As an alternate seminal approach [97], random Fourier features were used to develop approximate feature maps  $\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^Y$  for the RBF kernel such that  $\kappa(\mathbf{x}, \mathbf{x}') \approx \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$  [97]. More precisely, with high probability,

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^D, |\kappa(\mathbf{x}, \mathbf{x}') - \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle| \sim O\left(\sqrt{D/Y}\right), \quad (5.29)$$

implying that an  $\epsilon$  approximation guarantee requires  $Y \sim O(D/\epsilon^2)$ . The first set of random

feature maps for the RBF kernel were defined as follows using trigonometric functions:

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{Y}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}^1, \mathbf{x} \rangle + b^1) \\ \cos(\langle \boldsymbol{\omega}^2, \mathbf{x} \rangle + b^2) \\ \dots \\ \cos(\langle \boldsymbol{\omega}^Y, \mathbf{x} \rangle + b^Y) \end{bmatrix}, \text{ and } \Phi(\mathbf{x}) = \frac{1}{\sqrt{Y}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}^1, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^1, \mathbf{x} \rangle) \\ \cos(\langle \boldsymbol{\omega}^2, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^2, \mathbf{x} \rangle) \\ \dots \\ \cos(\langle \boldsymbol{\omega}^Y, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^Y, \mathbf{x} \rangle) \end{bmatrix}, \quad \begin{aligned} &\boldsymbol{\omega}^i \sim \mathcal{N}(0, \mathbf{I}_D), \\ &b^i \sim \mathcal{U}(0, 1), \\ &\forall i \in \llbracket Y \rrbracket. \end{aligned} \quad (5.30)$$

Note that the first set of random features produces a  $Y$ -dimensional feature map with  $Y$  random features while the second set produces a  $2Y$ -dimensional feature map using  $Y$  random features, but provides better approximation guarantees. The  $\mathcal{N}(0, \mathbf{I}_D)$  denotes the  $D$ -dimensional isotropic standard normal distribution, and the  $\mathcal{U}(0, 1)$  denotes the univariate uniform distribution over the range  $[0, 1]$ . Since then, various other random features have been developed for the RBF kernel [98; 99; 100; 101], and for various other kernels [102; 103]. Liu et al., 2021 [104] provide a comprehensive survey of random feature for kernel approximation.

In the context of Associative Memory, this allows us to approximate the energy function of the form in Eq. (5.27) as follows:

$$E_\beta(\mathbf{v}; \Xi) = -Q \left( \sum_{\mu \in \llbracket K \rrbracket} \kappa(\mathbf{v}, \boldsymbol{\xi}^\mu) \right) \approx -Q \left( \left\langle \Phi(\mathbf{v}), \underbrace{\sum_{\mu=1}^K \Phi(\boldsymbol{\xi}^\mu)}_{\triangleq \mathbf{T}} \right\rangle \right) = \tilde{E}_\beta(\mathbf{v}; \mathbf{T}), \quad (5.31)$$

where the computation of the energy  $E_\beta(\cdot; \Xi)$  (and its gradient) requires us to have access to all the stored patterns  $\Xi = \{\boldsymbol{\xi}^\mu, \mu \in \llbracket K \rrbracket\}$  of size  $KD$ , while the computation of the approximate energy  $\tilde{E}_\beta(\cdot; \mathbf{T})$  (and its gradient) only requires us to have access to the  $\mathbf{T} \triangleq \sum_{\mu=1}^K \Phi(\boldsymbol{\xi}^\mu)$  of size  $Y$  and the random features  $\{(\boldsymbol{\omega}^i, b^i), i \in \llbracket Y \rrbracket\}$  (which can be generated on the fly given the  $Y$  random seeds and thus do not need to be stored explicitly) [105]. We can now perform inference via gradient descent on this approximate energy, providing an unique (Dense) Associative Memory model that does not require the stored patterns  $\Xi$  for inference. It has been shown that the approximation in the energy translates to approximation in the inference — the inference  $f_\Xi(\mathbf{x})$  by minimizing the exact energy  $E_\beta(\cdot; \Xi)$  is approximated with the inference  $f_\mathbf{T}(\mathbf{x})$  by minimizing the approximate energy  $\tilde{E}_\beta(\cdot; \mathbf{T})$ . This approximation is affected by the following factors [105]:

- The approximation depends on the kernel approximation introduced by the random features with a factor of  $O(\sqrt{N/Y})$ , with larger number of random features  $Y$  improving the approximation.
- The approximation also depends on the initial energy  $E_\beta(\mathbf{x}; \Xi)$  of the input  $\mathbf{x}$  — the initial state for the energy descent — with larger initial energy leading to higher levels of approximation.
- The hyperparameter  $\eta$  which corresponds to the step-size (or learning rate) of the energy

gradient descent, with smaller  $\eta$  implying lower levels of approximation.

Of course, if we are already considering a kernel function  $\kappa$  in Eq. (5.27), which has an explicit feature map  $\phi$  (for example if  $\kappa(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$  or  $\kappa(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle)^2$ ), then we can directly use the explicit exact feature map to simplify the exact energy, and incur no approximation in the kernel function evaluation, and thus in the Associative Memory model inference.

### Notebook 5.1: Distributed Representation for Dense Associative Memory

In this notebook, we demonstrate how we utilize random features to disentangle the size of the Dense Associative Memory network from the number of memories to be stored. Given the standard log-sum-exp energy  $E_\beta(\cdot; \Xi)$ , corresponding to a model  $f_\Xi$  of size  $O(DK)$ , we demonstrate how we can use the trigonometric random features to develop an approximate energy  $\tilde{E}_\beta(\cdot; \mathbf{T})$  using a distributed representation  $\mathbf{T}$  of the memories  $\Xi = \{\xi^\mu, \mu \in \llbracket K \rrbracket\}$ , thus giving us a model  $f_{\mathbf{T}}$  of size  $O(Y)$ .

Checkout the notebook as a [blog post](#), a [colab notebook](#) or as a [raw .ipynb file](#).

$$\begin{aligned}
 \underbrace{E_\beta(\mathbf{v}; \Xi)}_{f_\Xi \sim O(DK) \text{ size}} &= -\log \sum_{\mu=1}^K \exp(-\beta/2 \|\mathbf{v} - \xi^\mu\|^2) \\
 &\approx -\log \sum_{\mu=1}^K \langle \Phi(\sqrt{\beta}\mathbf{v}), \Phi(\sqrt{\beta}\xi^\mu) \rangle \\
 &= -\log \left\langle \Phi(\sqrt{\beta}\mathbf{v}), \sum_{\mu=1}^K \Phi(\sqrt{\beta}\xi^\mu) \right\rangle \\
 &= -\log \langle \Phi(\sqrt{\beta}\mathbf{v}), \mathbf{T} \rangle = \underbrace{\tilde{E}_\beta(\mathbf{v}; \mathbf{T})}_{f_{\mathbf{T}} \sim O(Y) \text{ size}}
 \end{aligned}
 \quad
 \Phi(\mathbf{x}) = \frac{1}{\sqrt{Y}} \begin{bmatrix} \cos(\langle \omega^1, \mathbf{x} \rangle) \\ \sin(\langle \omega^1, \mathbf{x} \rangle) \\ \cos(\langle \omega^2, \mathbf{x} \rangle) \\ \sin(\langle \omega^2, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \omega^Y, \mathbf{x} \rangle) \\ \sin(\langle \omega^Y, \mathbf{x} \rangle) \end{bmatrix}$$

$\omega^i \sim \mathcal{N}(0, \mathbf{I}_D), \forall i \in \llbracket Y \rrbracket$

#### 5.4.2 Novel Energy Functions

As discussed earlier in Section (5.2), given an energy function, we can define a probability density through the Boltzmann distribution. Alternately, given a density function  $p : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ , we can define a energy function  $E(\mathbf{v}) \propto -\log p(\mathbf{v})$  through the same relationship.

Given a set of samples  $\Xi = \{\xi^\mu \sim p_{\text{data}}, \mu \in \llbracket K \rrbracket\}$  from an unknown distribution  $p_{\text{data}}$  over  $\mathcal{X} \subset \mathbb{R}^D$ , one way to define a density function  $\hat{p}$  is through kernel density estimation, where the goal is to devise a  $\hat{p}$  that closely approximates the unknown  $p_{\text{data}}$ . A kernel density estimate or KDE at any point  $\mathbf{v} \in \mathcal{X}$  is defined as:

$$\hat{p}_h(\mathbf{v}; \Xi) = \frac{1}{Kh} \sum_{\mu=1}^K \kappa\left(\frac{\mathbf{v} - \xi^\mu}{h}\right), \quad (5.32)$$

where  $\kappa : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  is the kernel function, and  $h > 0$  is the kernel bandwidth. For this to be a valid density, the kernel function needs to satisfy the following conditions:

- Symmetry:  $\kappa(\mathbf{x}) = \kappa(-\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}$

- Nonnegativity:  $\kappa(\mathbf{x}) \geq 0 \quad \forall \mathbf{x} \in \mathcal{X}$
- Normalization:  $\int_{\mathbf{x}} \kappa(\mathbf{x}) d\mathbf{x} = 1$ .

For multivariate data (that is  $D > 1$ ), the kernel  $\kappa$  has been defined both as  $\kappa(\mathbf{x}) = c\kappa_1(\|\mathbf{x}\|)$  or  $\kappa(\mathbf{x}) = c' \prod_{i=1}^D \kappa_1(|x_i|)$ , where  $\kappa_1 : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  is an univariate kernel function,  $x_i$  denotes the  $i$ -th coordinate of  $\mathbf{x}$  for any  $\mathbf{x} \in \mathbb{R}^D$ , and  $c, c'$  are positive constants ensuring that normalization condition for  $\kappa$  is satisfied. Note that, with the RBF kernel (which becomes the Gaussian kernel with proper scaling for normalization) with  $\kappa_1(z) = \exp(-|z|^2)$ ,  $z \in \mathbb{R}$ ,  $\kappa_1(\|\mathbf{x}\|) = \prod_i \kappa_1(|x_i|) = \exp(-\|\mathbf{x}\|^2)$ . We will consider the univariate case from hereon for the ease of exposition with  $D = 1$ , where  $\kappa : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  satisfying the aforementioned symmetry, nonnegativity and normalization conditions.

For the purpose of KDE, the scale of the kernel function is not unique. That is, for a given  $\kappa(\cdot)$ , we can define  $\tilde{\kappa}(\cdot) = b^{-1}\kappa(\cdot/b)$ , for some  $b > 0$ . Then, one obtains the same KDE by rescaling the choice of  $h$ . Hence, the shape of the kernel function plays a more important role in determining the choice of the kernel. We now introduce two parameters associated with the kernel – the *scale*  $\mu_\kappa$  and the *regularity*  $\sigma_\kappa$  defined as:

$$\mu_\kappa \triangleq \int_{\mathbb{R}} x^2 \kappa(x) dx, \quad \sigma_\kappa \triangleq \int_{\mathbb{R}} (\kappa(x))^2 dx \quad (5.33)$$

The quality or generalization of KDE depends on these two properties of the kernel. The generalization error of  $\hat{p}_h(\cdot; \Xi)$  is measured by the *Mean Integrated Squared Error* or MISE, and is given by

$$\text{MISE}(h) = \mathbb{E} \left[ \int_v (\hat{p}_h(v; \Xi) - p_{\text{data}}(v))^2 dv \right], \quad (5.34)$$

where the expectation is over the  $K$  random samples  $\Xi$  from  $p_{\text{data}}$ .

Assuming that the ground-truth density  $p_{\text{data}}$  is twice continuously differentiable, a second-order Taylor expansion gives the leading terms of the  $\text{MISE}(h)$ , which decomposes into squared bias and variance terms [106, Section 2.5]:

$$\text{MISE}(h) \approx \underbrace{\frac{\mu_\kappa^2}{4} h^4 \int_v |p''_{\text{data}}(v)|^2 dv}_{\text{bias-squared term}} + \underbrace{\frac{\sigma_\kappa}{Kh}}_{\text{variance}}. \quad (5.35)$$

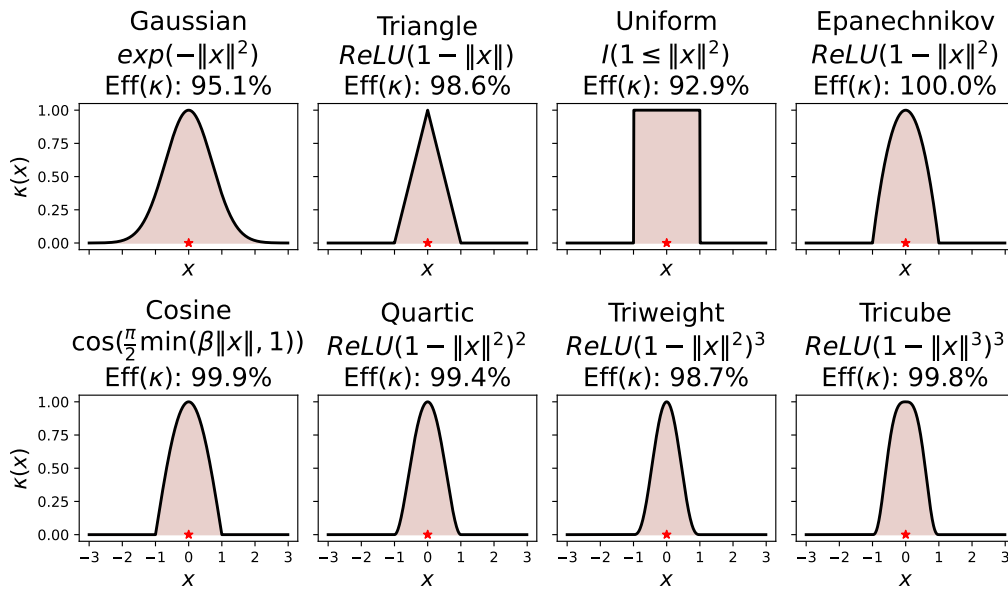
Thus, reducing the bandwidth  $h$  decreases bias but increase variance, and vice verse for increasing  $h$ , thereby highlighting the bias-variance tradeoff. Balancing the bias-squared and variance terms, we can have kernel-specific optimal choice  $h_\kappa^*$  for the bandwidth

$$h_\kappa^* \triangleq \left( \frac{\sigma_\kappa}{K\mu_\kappa^2} \frac{4}{\int_v |p''_{\text{data}}(v)|^2 dv} \right)^{1/5}. \quad (5.36)$$

Plugging this into Eq. (5.35) gives us the best possible MISE:

$$\text{MISE}(h_\kappa^*) \approx \frac{5}{4} \left( \frac{\sqrt{\mu_\kappa} \sigma_\kappa \int_v |p''_{\text{data}}(v)|^2 dv}{K} \right)^{4/5}, \quad (5.37)$$

where the choice of the kernel function  $\kappa$  affects the MISE through its scale  $\mu_\kappa$  and regularity  $\sigma_\kappa$ . Thus, it is intuitive to select the kernel function  $\kappa$  based on the optimal  $\text{MISE}(h_\kappa^*)$ . As discussed above, the scale of the kernel function is non-unique, and can be fixed to  $\mu_\kappa = 1$  by appropriately scaling the kernel function. Hence, the kernel with the smallest regularity  $\sigma_\kappa$ , subject to  $\mu_\kappa = 1$  (without loss of generality), over the class of normalized, symmetric, and positive kernels is most desirable. This is a well-studied problem [107; 108] [106, Section 2.7], and the Epanechnikov kernel  $\kappa_{\text{epan}}(z) = \max\{1 - |z|^2, 0\} = \text{ReLU}(1 - |z|^2)$  achieves the optimal  $\text{MISE}(h_\kappa^*)$ . The quantity,  $\text{Eff}(\kappa) \triangleq \sigma_\kappa / \sigma_{\kappa_{\text{epan}}}$  is known as the *efficiency* of any kernel relative to the Epanechnikov kernel. Various kernels with varying levels of efficiencies have been developed, and we present a representative subset of these kernel functions in Fig. (5.7). Similar analysis and guarantees can be established for multivariate KDE.



**Figure 5.7: Different kernels used for Kernel Density Estimation.** We show the shapes, expression and KDE efficiency relative to the Epanechnikov kernel (*higher is better*) for 8 kernels. The center of each kernel is marked with a red  $\star$ . To highlight the shape of the kernel, we have removed any scaling in the kernel expression. Note that all above kernels except Gaussian have finite support. The Epanechnikov kernel has the highest efficiency (100%). The Gaussian kernel is extremely popular, and it is more efficient (95.1%) than the Uniform kernel (92.9%). However, there are various other kernels (such as the Triangle kernel) with better efficiency. This image is replicated from Hoover et al., 2025 [109].

The rich literature of KDE and its suite of well-studied kernel functions opens up the path to the development of various energy functions for Associative Memory networks — one for each kernel function — which have not been considered previously. As a natural first choice, one can select the optimal Epanechnikov kernel, which leads to the following novel energy function [109]:

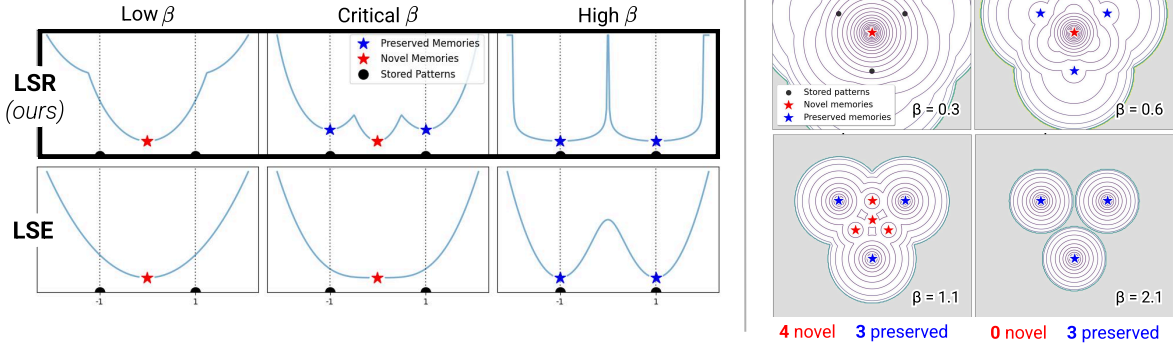
$$E_\beta(\mathbf{v}; \Xi) = -\log \sum_{\mu=1}^K \text{ReLU}\left(1 - \beta/2 \|\mathbf{v} - \xi^\mu\|^2\right). \quad (5.38)$$

This makes use of a *shifted-ReLU* operation, and thus is termed the log-sum-ReLU or LSR energy.

This can be contrasted with the popular LSE or log-sum-exp energy shown in Eq. (5.13), where we replace the exponential separation function  $F(z) = \exp(z)$  with the shifted-ReLU separation function  $F(z) = \text{ReLU}(1 + z)$  with the negative squared Euclidean distance based similarity function  $S[\mathbf{x}, \mathbf{x}'] = -1/2 \|\mathbf{x} - \mathbf{x}'\|^2$ .

**LSR preserves** memories while creating **novel** ones.

*LSE can do only one or the other.*



**Figure 5.8: Emergence of novel energy local minima.** LSR energy can create more memories than there are stored patterns under critical regimes of  $\beta$ . **Left:** 1D LSR vs LSE energy landscape. Note that LSE is never capable of having more local minima than the number of stored patterns. **Right:** 2D LSR energy landscape, where increasing  $\beta$  creates novel local minima where basins intersect. Unsupported regions are shaded gray. This image is replicated from Hoover et al., 2025 [109].

This novel energy function has various desirable properties:

#### Exact single-step retrieval

For an Associative Memory network with the LSR energy and appropriate hyperparameters, it is possible to have exact retrieval of stored patterns in a single energy gradient step. This is in contrast to LSE where only approximate retrieval is possible unless the inverse-temperature  $\beta \rightarrow \infty$ .

#### Exponential memory capacity without exponential separation function

This Associative Memory network equipped with the LSR energy has exponential memory capacity — that is, the number of stored patterns that are retrievable is  $O(\exp(D))$ . This is similar to the LSE energy.

#### Generation of a multitude of novel memories

Finally, the LSE energy can introduce numerous novel energy local minima to the energy landscape, while also maintaining local minima around the stored patterns, enabling simultaneous retrieval of stored patterns and retrieval of novel memory, providing a path to data generation in Associative Memory networks with energy descent. This phenomena is visualized for data in one and two dimensions in Fig. (5.8), and has been utilized to create

novel samples from an approximation of the underlying data distribution  $p_{\text{data}}$ . Such a phenomena has not previously been seen in literature.

However, this novel LSR energy can pose certain novel challenges:

### Regions of infinite energy

For a given configuration of an Associative Memory network, with LSR energy, there exists  $\mathbf{v} \in \mathcal{X}$  such that  $E_{\beta}(\mathbf{v}; \Xi) = \infty$  given the finite support of the Epanechnikov kernel. This is visualized in Fig. (5.8) (Right) as the gray shaded region for two dimensional data.

## Chapter 6

# Conclusion

In this tutorial we have covered recent advances in energy-based Associative Memory, including information storage capacity calculations (Chapter 2), relationship to transformers (Chapter 3) and diffusion models (Chapter 4), and connection to non-neural network machine learning (Chapter 5) and other ideas. In the past few years, Associative Memory became an active area of research with many lines of exploration coexisting and branching into several disciplines. Since this tutorial was prepared for ICML audience, we focused mostly on explaining the core ideas with only minimal derivations necessary to understand those ideas. We also prepared coding notebooks to help AI practitioners gain hands-on experience with AM basics. Inevitably, with this strategy in mind, many important and exciting aspects of AMs remained outside the scope of this tutorial. For instance, we did not discuss the biological implementations of DenseAMs [30; 110; 111; 35; 112; 113; 114]. Several valuable trends in AM-inspired statistical physics [115; 116; 117; 118; 28; 119; 120; 121; 122] have also only been briefly mentioned. Memory augmentation of large language models (LLMs) [123; 124; 125; 126] is becoming an active area of research with clever ideas on how memory models can be used synergistically with feed-forward architectures. There are exciting ideas around novel neural architectures inspired by AMs [29; 33; 127; 128; 26; 129; 130], and domain specific applications [131; 132] that have not been covered with sufficient detail either. Quantum DenseAMs is an emerging topic [133]. Neuromorphic hardware based on DenseAMs [134] is becoming a promising area of research too. We expect these trends to grow and new trends to appear. We hope that this introductory tutorial may provide an entry point for new researchers in this exciting field.

# Bibliography

- [1] Endel Tulving. Episodic and semantic memory. *Organization of memory*, 1972.
- [2] Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- [3] Frank Rosenblatt. Principles of neurodynamics. *Perceptrons and the theory of brain mechanisms*, 1962.
- [4] The New York Times. New navy device learns by doing; psychologist shows embryo of computer designed to read and grow wiser. *The New York Times*, 1958.
- [5] Minsky Marvin and A Papert Seymour. Perceptrons. *Cambridge, MA: MIT Press*, 6(318-362):7, 1969.
- [6] The Royal Swedish Academy of Sciences. Nobel prize in physics 2024. *Nobel Prize Outreach*, 2024.
- [7] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- [8] Stephen G. Brush. History of the lenz-ising model. *Rev. Mod. Phys.*, 39:883–893, 10 1967.
- [9] James A Anderson. A memory storage model utilizing spatial correlation functions. *Kybernetik*, 5(3):113–119, 1968.
- [10] David J Willshaw, O Peter Buneman, and Hugh Christopher Longuet-Higgins. Non-holographic associative memory. *Nature*, 222(5197):960–962, 1969.
- [11] S-I Amari. Learning patterns and pattern sequences by self-organizing nets of threshold elements. *IEEE Transactions on computers*, 100(11):1197–1206, 1972.
- [12] S-I Amari. Neural theory of association and concept-formation. *Biological cybernetics*, 26(3):175–185, 1977.
- [13] Michael A Cohen and Stephen Grossberg. Absolute stability of global pattern formation and parallel memory storage by competitive neural networks. *IEEE transactions on systems, man, and cybernetics*, 5:815–826, 1983.
- [14] John Hopfield. Neurons With Graded Response Have Collective Computational Properties Like Those of Two-State Neurons. *Proceedings of the National Academy of Sciences of the United States of America*, 81:3088–92, June 1984.

- [15] Daniel J Amit, Hanoch Gutfreund, and Haim Sompolinsky. Storing infinite numbers of patterns in a spin-glass model of neural networks. *Physical Review Letters*, 55(14):1530, 1985.
- [16] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [17] Paul J. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1(4):339–356, 1988.
- [18] Donald O. Hebb. (1949) donald o. hebb, the organization of behavior, new york: Wiley, introduction and chapter 4, "the first stage of perception: growth of the assembly," pp. xi-xix, 60-78. In *Neurocomputing, Volume 1: Foundations of Research*. The MIT Press, 04 1949.
- [19] Geoffrey Hinton and Terrence Sejnowski. Optimal perceptual inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 448–453, 01 1983.
- [20] John J. Hopfield, David I. Feinstein, and Richard G. Palmer. ‘unlearning’ has a stabilizing effect in collective memories. *Nature*, 304:158–159, 1983.
- [21] Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29, 2016.
- [22] Kuniyiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36:193–202, 1980.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [24] Dmitry Krotov. A new frontier for hopfield networks. *Nature Reviews Physics*, 5(7):366–367, 2023.
- [25] Dmitry Krotov and John Hopfield. Dense associative memory is robust to adversarial inputs. *Neural computation*, 30(12):3151–3167, 2018.
- [26] Hamza Chaudhry, Jacob Zavatone-Veth, Dmitry Krotov, and Cengiz Pehlevan. Long sequence hopfield memory. *Advances in Neural Information Processing Systems*, 36:54300–54340, 2023.
- [27] 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, 2017.
- [28] Carlo Lucibello and Marc Mézard. Exponential capacity of dense associative memories. *Physical Review Letters*, 132(7):077301, 2024.
- [29] 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.
- [30] Dmitry Krotov and John J Hopfield. Large associative memory problem in neurobiology and machine learning. In *International Conference on Learning Representations*, 2021.
- [31] Beren Millidge, Tommaso Salvatori, Yuhang Song, Thomas Lukasiewicz, and Rafal Bogacz. Universal hopfield networks: A general framework for single-shot associative memory models. In *International Conference on Machine Learning*, pages 15561–15583. PMLR, 2022.
- [32] Thomas F Burns and Tomoki Fukai. Simplicial hopfield networks. In *The Eleventh International Conference on Learning Representations*, 2022.
- [33] Dmitry Krotov. Hierarchical associative memory, 2021.
- [34] Benjamin Hoover, Yuchen Liang, Bao Pham, Rameswar Panda, Hendrik Strobelt, Duen Horng Chau, Mohammed Zaki, and Dmitry Krotov. Energy transformer. *Advances in Neural Information Processing Systems*, 36, 2024.
- [35] Leo Kozachkov, Jean-Jacques Slotine, and Dmitry Krotov. Neuron–astrocyte associative memory. *Proceedings of the National Academy of Sciences*, 122(21):e2417788122, 2025.
- [36] Benjamin Hoover, Duen Horng Chau, Hendrik Strobelt, and Dmitry Krotov. A universal abstraction for hierarchical hopfield networks. *The Symbiosis of Deep Learning and Differential Equations II*, 2022.
- [37] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [38] Fei Tang and Michael Kopp. A remark on a paper of krotov and hopfield. *arXiv preprint arXiv:2105.15034*, 2021.
- [39] Roger Brown and James Kulik. Flashbulb memories. *Cognition*, 5(1):73–99, 1977.
- [40] Marigold Linton. I remember it well. *Psychology Today*, 13(2):81, 1979.
- [41] Elizabeth F Loftus. *Memory*. Rowman & Littlefield Publishers, 1988.
- [42] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 7 2015. PMLR.
- [43] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [44] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.

- [45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [46] Benjamin Hoover, Hendrik Strobelt, Dmitry Krotov, Judy Hoffman, Zolt Kira, and Duen Horng Chau. Memory in Plain Sight: Surveying the Uncanny Resemblances of Associative Memories and Diffusion Models, 2023.
- [47] Luca Ambrogioni. In Search of Dispersed Memories: Generative Diffusion Models Are Associative Memory Networks. *Entropy*, 26(5), 2024.
- [48] Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023.
- [49] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- [50] Bao Pham, Gabriel Raya, Matteo Negri, Mohammed J Zaki, Luca Ambrogioni, and Dmitry Krotov. Memorization to generalization: Emergence of diffusion models from associative memory. *arXiv preprint arXiv:2505.21777*, 2025.
- [51] John Hertz, Anders Krogh, and Richard G Palmer. *Introduction to the theory of neural computation*. Addison Wesley Longman, 1991.
- [52] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*, 2020.
- [53] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- [54] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pages 21450–21474, 2023.
- [55] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.
- [56] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- [57] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent

- diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [58] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. *openai.com*, 2024.
  - [59] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Understanding and mitigating copying in diffusion models. *Advances in Neural Information Processing Systems*, 36:47783–47803, 2023.
  - [60] Bishwajit Saha, Dmitry Krotov, Mohammed J Zaki, and Parikshit Ram. End-to-end differentiable clustering with associative memories. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 29649–29670. PMLR, 23–29 Jul 2023.
  - [61] C Cortes, A Krogh, and JA Hertz. Hierarchical associative networks. *Journal of Physics A: Mathematical and General*, 20(13):4449, 1987.
  - [62] Hanoch Gutfreund. Neural networks with hierarchically correlated patterns. *Physical Review A*, 37(2):570, 1988.
  - [63] A Krogh and JA Hertz. Mean-field analysis of hierarchical associative networks with ‘magnetisation’. *Journal of Physics A: Mathematical and General*, 21(9):2211, 1988.
  - [64] I Kanter and Haim Sompolinsky. Associative recall of memory without errors. *Physical Review A*, 35(1):380, 1987.
  - [65] JL Van Hemmen. Hebbian learning, its correlation catastrophe, and unlearning. *Network: Computation in Neural Systems*, 8(3):V1, 1997.
  - [66] Aditya Cowsik and Adithya Sriram. Dense hopfield networks with hierarchical memories. In *New Frontiers in Associative Memories*, 2025.
  - [67] Izrail Solomonovich Gradshteyn and Iosif Moiseevich Ryzhik. *Table of integrals, series, and products*. Academic press, 2014.
  - [68] Casey Meehan, Kamalika Chaudhuri, and Sanjoy Dasgupta. A non-parametric test to detect data-copying in generative models. In *International Conference on Artificial Intelligence and Statistics*, 2020.
  - [69] Gerrit J.J. Van den Burg and Chris Williams. On memorization in probabilistic deep generative models. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
  - [70] TaeHo Yoon, Joo Young Choi, Sehyun Kwon, and Ernest K Ryu. Diffusion probabilistic models generalize when they fail to memorize. In *ICML 2023 Workshop on Structured Probabilistic Inference Generative Modeling*, 2023.
  - [71] Xiangming Gu, Chao Du, Tianyu Pang, Chongxuan Li, Min Lin, and Ye Wang. On memorization in diffusion models. *arXiv preprint arXiv:2310.02664*, 2023.

- [72] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *Proceedings of the 32nd USENIX Conference on Security Symposium*, SEC '23, USA, 2023. USENIX Association.
- [73] Beatrice Achilli, Enrico Ventura, Gianluigi Silvestri, Bao Pham, Gabriel Raya, Dmitry Krotov, Carlo Lucibello, and Luca Ambrogioni. Losing dimensions: Geometric memorization in generative diffusion. *arXiv preprint arXiv:2410.08727*, 2024.
- [74] Giulio Biroli, Tony Bonnaire, Valentin de Bortoli, and Marc Mézard. Dynamical regimes of diffusion models. *Nature Communications*, 15(1), November 2024.
- [75] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Mach. Learn.*, 20(3):273–297, September 1995.
- [76] J. A. Nelder and R. W. M. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)*, 135(3):370–384, 1972.
- [77] J. R. Quinlan. Induction of decision trees. *Machine Learning*, 1:81–106, 1986.
- [78] Vladimir Vapnik. *Estimation of dependences based on empirical data*. Springer Science & Business Media, 2006.
- [79] Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31. Springer Science & Business Media, 2013.
- [80] Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. In *Advances in neural information processing systems*, pages 161–168, 2008.
- [81] Parikshit Ram, Alexander G Gray, Horst C Samulowitz, and Gregory Bramble. Toward theoretical guidance for two common questions in practical cross-validation based hyperparameter selection. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pages 802–810. SIAM, 2023.
- [82] Jerry Yao-Chieh Hu, Bo-Yu Chen, Dennis Wu, Feng Ruan, and Han Liu. Nonparametric modern hopfield models. *arXiv preprint arXiv:2404.03900*, 2024.
- [83] Dennis Wu, Jerry Yao-Chieh Hu, Teng-Yun Hsiao, and Han Liu. Uniform memory retrieval with larger capacity for modern hopfield models. *arXiv preprint arXiv:2404.03900*, 2024.
- [84] 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. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [85] Sanjoy Dasgupta. The hardness of k-means clustering. *UCSD Technical Report*, 2008.
- [86] Stuart Lloyd. Least squares quantization in PCM. *IEEE transactions on information theory*, 28(2):129–137, 1982.

- [87] J. C. Dunn. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics*, 3:32–57, 1974.
- [88] James C Bezdek. *Pattern recognition with fuzzy objective function algorithms*. Springer Science & Business Media, 2013.
- [89] Bishwajit Saha, Dmitry Krotov, Mohammed J Zaki, and Parikshit Ram. Deep clustering with associative memories. In *NeurIPS Workshop on Machine Learning and Compression*, 2024.
- [90] Erxue Min, Xifeng Guo, Qiang Liu, Gen Zhang, Jianjing Cui, and Jun Long. A survey of clustering with deep learning: From the perspective of network architecture. *IEEE Access*, 6:39501–39514, 2018.
- [91] Yazhou Ren, Jingyu Pu, Zhimeng Yang, Jie Xu, Guofeng Li, Xiaorong Pu, S Yu Philip, and Lifang He. Deep clustering: A comprehensive survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [92] Sheng Zhou, Hongjia Xu, Zhuonan Zheng, Jiawei Chen, Zhao Li, Jiajun Bu, Jia Wu, Xin Wang, Wenwu Zhu, and Martin Ester. A comprehensive survey on deep clustering: Taxonomy, challenges, and future directions. *ACM Computing Surveys*, 57(3), November 2024.
- [93] Léon Bottou. *Large-scale kernel machines*. MIT press, 2007.
- [94] Parikshit Ram, Dongryeol Lee, William March, and Alexander Gray. Linear-time algorithms for pairwise statistical problems. *Advances in Neural Information Processing Systems*, 22, 2009.
- [95] Ryan Curtin, William March, Parikshit Ram, David Anderson, Alexander Gray, and Charles Isbell. Tree-independent dual-tree algorithms. In *International Conference on Machine Learning*, pages 1435–1443. PMLR, 2013.
- [96] Ryan R Curtin, Dongryeol Lee, William B March, and Parikshit Ram. Plug-and-play dual-tree algorithm runtime analysis. *J. Mach. Learn. Res.*, 16:3269–3297, 2015.
- [97] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 2007.
- [98] Felix Xinnan X Yu, Ananda Theertha Suresh, Krzysztof M Choromanski, Daniel N Holtmann-Rice, and Sanjiv Kumar. Orthogonal random features. *Advances in neural information processing systems*, 29, 2016.
- [99] Krzysztof M Choromanski, Mark Rowland, and Adrian Weller. The unreasonable effectiveness of structured random orthogonal embeddings. *Advances in neural information processing systems*, 30, 2017.
- [100] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *Proceedings of ICLR*, 2020.

- [101] Valerii Likhoshesterov, Krzysztof Marcin Choromanski, Kumar Avinava Dubey, Frederick Liu, Tamas Sarlos, and Adrian Weller. Dense-exponential random features: Sharp positive estimators of the gaussian kernel. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [102] Purushottam Kar and Harish Karnick. Random feature maps for dot product kernels. In *Artificial intelligence and statistics*, pages 583–591. PMLR, 2012.
- [103] Raffay Hamid, Ying Xiao, Alex Gittens, and Dennis DeCoste. Compact random feature maps. In *International conference on machine learning*, pages 19–27. PMLR, 2014.
- [104] Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan AK Suykens. Random features for kernel approximation: A survey on algorithms, theory, and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):7128–7148, 2021.
- [105] Benjamin Hoover, Duen Horng Chau, Hendrik Strobelt, Parikshit Ram, and Dmitry Krotov. Dense associative memory through the lens of random features. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [106] Matt P Wand and M Chris Jones. *Kernel smoothing*. CRC press, 1994.
- [107] Vassiliy A Epanechnikov. Non-parametric estimation of a multivariate probability density. *Theory of Probability & Its Applications*, 14(1):153–158, 1969.
- [108] Hans-Georg Müller. Smooth optimum kernel estimators of densities, regression curves and modes. *The Annals of Statistics*, pages 766–774, 1984.
- [109] Benjamin Hoover, Krishna Balasubramanian, Dmitry Krotov, and Parikshit Ram. Dense associative memory with epanechnikov energy. In *New Frontiers in Associative Memories*, 2025.
- [110] Mallory A Snow and Jeff Orchard. Biological softmax: Demonstrated in modern hopfield networks. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 44, 2022.
- [111] Danil Tyulmankov, Kim Stachenfeld, Dmitry Krotov, and Larry Abbott. Memorization and consolidation in associative memory networks. In *Associative Memory {&E} Hopfield Networks in 2023*, 2023.
- [112] Sarthak Chandra, Sugandha Sharma, Rishidev Chaudhuri, and Ila Fiete. Episodic and associative memory from spatial scaffolds in the hippocampus. *Nature*, pages 1–13, 2025.
- [113] Mohadeseh Shafiei Kafraji, Dmitry Krotov, Brendan A Bicknell, and Peter E Latham. A biologically plausible associative memory network. In *New Frontiers in Associative Memories*, 2025.
- [114] Kaining Zhang and Gaia Tavoni. Maximizing memory capacity in heterogeneous networks. *PRX Life*, 3(2):023016, 2025.
- [115] Elena Agliari, Francesco Alemanno, Adriano Barra, Martino Centonze, and Alberto Fachechi.

- Neural networks with a redundant representation: Detecting the undetectable. *Physical review letters*, 124(2):028301, 2020.
- [116] Elena Agliari and Giordano De Marzo. Tolerance versus synaptic noise in dense associative memories. *The European Physical Journal Plus*, 135(11):1–22, 2020.
- [117] Linda Albanese, Francesco Alemanno, Andrea Alessandrelli, and Adriano Barra. Replica symmetry breaking in dense hebbian neural networks. *Journal of Statistical Physics*, 189(2):24, 2022.
- [118] Elena Agliari, Linda Albanese, Francesco Alemanno, Andrea Alessandrelli, Adriano Barra, Fosca Giannotti, Daniele Lotito, and Dino Pedreschi. Dense hebbian neural networks: a replica symmetric picture of supervised learning. *Physica A: Statistical Mechanics and its Applications*, 626:129076, 2023.
- [119] Robin Thériault and Daniele Tantari. Dense hopfield networks in the teacher-student setting. *SciPost Physics*, 17(2):040, 2024.
- [120] David G Clark. Transient dynamics of associative memory models. *arXiv preprint arXiv:2506.05303*, 2025.
- [121] Kazushi Mimura, Jun’ichi Takeuchi, Yuto Sumikawa, Yoshiyuki Kabashima, and Anthony CC Coolen. Dynamical properties of dense associative memory. *arXiv preprint arXiv:2506.00851*, 2025.
- [122] Flavio Nicoletti, Francesco D’Amico, and Matteo Negri. Statistical mechanics of vector hopfield network near and above saturation. *arXiv preprint arXiv:2507.02586*, 2025.
- [123] Mikhail S Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V Sapunov. Memory transformer. *arXiv preprint arXiv:2006.11527*, 2020.
- [124] Zexue He, Leonid Karlinsky, Donghyun Kim, Julian McAuley, Dmitry Krotov, and Rogerio Feris. Camelot: Towards large language models with training-free consolidated associative memory. *arXiv preprint arXiv:2402.13449*, 2024.
- [125] Ivan Rodkin, Yuri Kuratov, Aydar Bulatov, and Mikhail Burtsev. Associative recurrent memory transformer. *arXiv preprint arXiv:2407.04841*, 2024.
- [126] Yu Wang, Dmitry Krotov, Yuanzhe Hu, Yifan Gao, Wangchunshu Zhou, Julian McAuley, Dan Gutfreund, Rogerio Feris, and Zexue He. M+: Extending memoryllm with scalable long-term memory. *arXiv preprint arXiv:2502.00592*, 2025.
- [127] Andreas Füst, Elisabeth Rumetshofer, Johannes Lehner, Viet T Tran, Fei Tang, Hubert Ramsauer, David Kreil, Michael Kopp, Günter Klambauer, Angela Bitto, et al. Cloob: Modern hopfield networks with infoloob outperform clip. *Advances in neural information processing systems*, 35:20450–20468, 2022.
- [128] Yuchen Liang, Dmitry Krotov, and Mohammed J Zaki. Modern hopfield networks for graph embedding. *Frontiers in big Data*, 5:1044709, 2022.

- [129] Ryo Karakida, Toshihiro Ota, and Masato Taki. Hierarchical associative memory, parallelized mlp-mixer, and symmetry breaking. *arXiv preprint arXiv:2406.12220*, 2024.
- [130] Xueyan Niu, Bo Bai, Lei Deng, and Wei Han. Beyond scaling laws: Understanding transformer performance with associative memory. *arXiv preprint arXiv:2405.08707*, 2024.
- [131] Michael Widrich, Bernhard Schäfl, Milena Pavlović, Hubert Ramsauer, Lukas Gruber, Markus Holzleitner, Johannes Brandstetter, Geir Kjetil Sandve, Victor Greiff, Sepp Hochreiter, et al. Modern hopfield networks and attention for immune repertoire classification. *Advances in neural information processing systems*, 33:18832–18845, 2020.
- [132] Qian Zhang, Dmitry Krotov, and George Em Karniadakis. Operator learning for reconstructing flow fields from sparse measurements: an energy transformer approach. *arXiv preprint arXiv:2501.08339*, 2025.
- [133] Takeshi Kimura and Kohtaro Kato. Analysis of discrete modern hopfield networks in open quantum system. *arXiv preprint arXiv:2411.02883*, 2024.
- [134] Khalid Musa, Santosh Kumar, Michael Katidis, and Yu-Ping Huang. Dense associative memory in a nonlinear optical hopfield neural network. *arXiv preprint arXiv:2506.07849*, 2025.