



Modern Methods in Associative Memory

AM and Broader ML|AI

Parikshit Ram
IBM Research

**How to think about
Associative Memory
Networks as a
Machine Learning
Model?**



**How to use them for
standard Machine
Learning tasks?**

Chapter 5

Associative Memory:
A Machine Learning Model

Associative Memory Networks

$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$

$$\Xi = [\xi^1, \xi^2, \dots, \xi^K], \xi^\mu \in \mathbb{R}^D$$

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

Model is parameterized with **stored patterns**

Energy function defines the **"model architecture"**

Monotonic
function

Separation
function

Similarity b/w
state & memory

Associative Memory Networks

$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$

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

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

$$f_{\Xi}(\mathbf{x}) \triangleq \mathbf{v}^{(T)}.$$

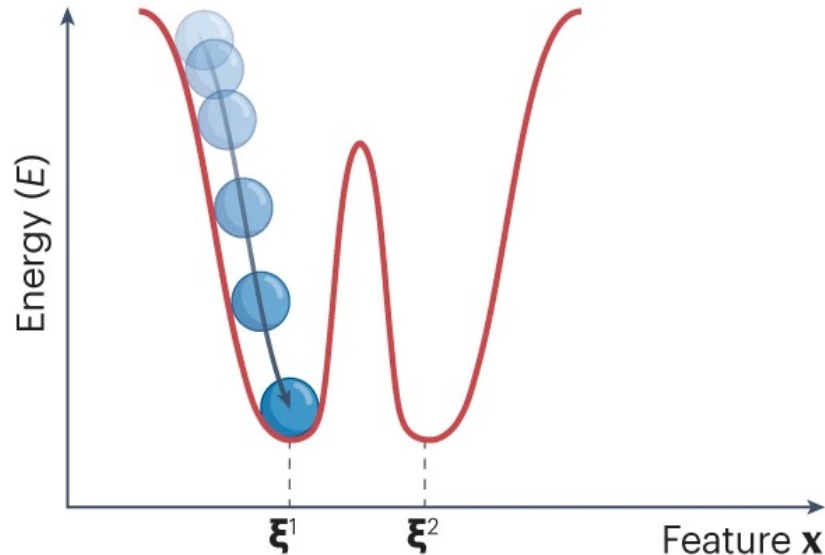
$$\Xi = [\xi^1, \xi^2, \dots, \xi^K], \xi^{\mu} \in \mathbb{R}^D$$

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

Inference (or the forward pass) with a T -layer AM network is equivalent to T **energy descent steps** using the energy gradient.

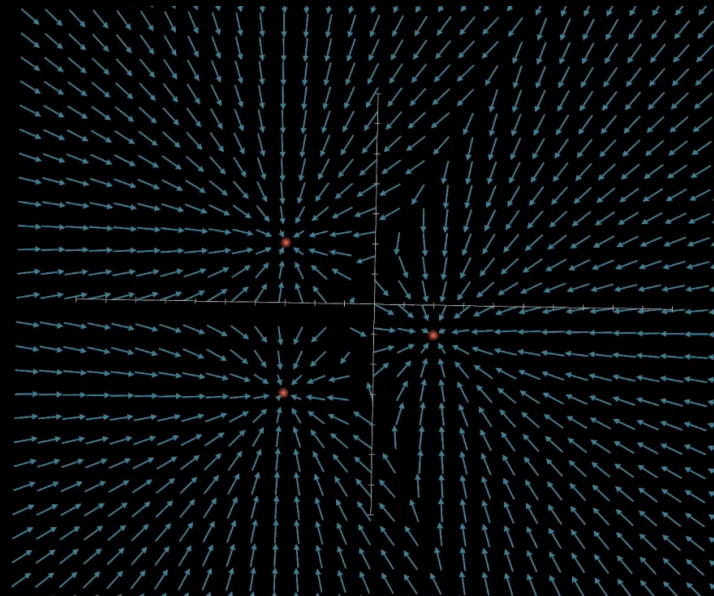
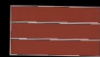
Associative Memory Networks

- The **energy of any state** is inversely related to the **probability** of the state
- Inference via **energy minimization** – or **likelihood maximization**



Associative Memory Networks

- More memories lead to more local minima of the energy *up to a point*
- More memories – more parameters – potentially more expressive model



Memory Capacity

The number of memories that can form distinct local minima of the energy

Classical Associative Memory

$$F(z) = z^2 \Rightarrow K^{\max} \sim O(D)$$

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

Dense Associative Memory

$$F(z) = z^N \Rightarrow K^{\max} \sim O(D^{N-1})$$

$$F(z) = \exp(z) \Rightarrow K^{\max} \sim O(\exp(D))$$

More memory capacity --
more model expressivity

But increased computational overhead!

Energy is a Kernel Sum

$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D \quad \Xi = [\xi^1, \xi^2, \dots, \xi^K], \xi^\mu \in \mathbb{R}^D$$

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

Kernel sum

$$F(\beta S[\mathbf{x}, \mathbf{x}']) \triangleq \kappa(\mathbf{x}, \mathbf{x}')$$

AMs and Kernel Machines

- Both **nonparametric**
- Kernel machines
 - Inference with a **single kernel sum** (usually)
- Associative Memory networks
 - AM can be **parametric**
 - Kernel sum computes energy & inference via energy descent
 - Single inference (usually) needs **multiple kernel sums**
 - Need ability to **differentiate through the kernel sums**

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

Kernel sum

AMs and Kernel Machines

Kernel machines

$$x \longrightarrow \sum_{\mu} \kappa(x, \xi^{\mu}) \longrightarrow \text{output}$$

$$E_{\beta}(\mathbf{v}; \Xi) = -Q \left[\sum_{\mu=1}^K \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu}) \right]$$

Kernel sum

Associative Memory networks

$$x \longrightarrow v^{(0)} \longrightarrow \begin{array}{c} v^{(t+1)} = v^{(t)} - \eta \nabla_v Q[s] \\ \text{\scriptsize $s = \sum_{\mu} \kappa(v^{(t)}, \xi^{\mu})$} \end{array} \longrightarrow \text{output}$$

AMs and Kernel Machines

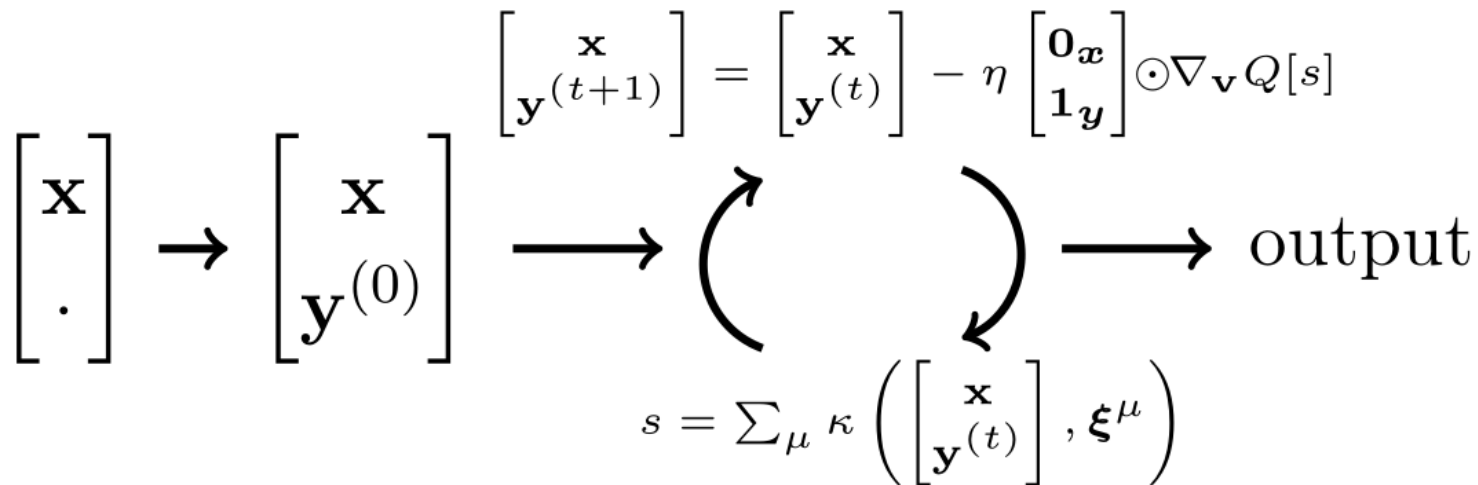
Kernel machines

$$\underbrace{\sum_{\mu} \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu})}_{\text{Density } p_{\mathbf{v}}(\mathbf{v})} \quad \longrightarrow \quad \underbrace{\sum_{\mu} \mathbf{w}_{\mu} \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu})}_{\int_{\mathbf{y}} \mathbf{y} p_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{y}} \quad \longrightarrow \quad \underbrace{\frac{\sum_{\mu} \mathbf{w}_{\mu} \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu})}{\sum_{\mu} \kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu})}}_{\mathbb{E}[p_{\mathbf{y}}(\mathbf{y}|\mathbf{x})] = \frac{\int_{\mathbf{y}} \mathbf{y} p_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{y}}{p_{\mathbf{x}}(\mathbf{x})}}$$

Kernel machines can compute the **expectation of the conditional distribution**.

AMs and Kernel Machines

Associative Memory networks



AM networks can find the **modes of the conditional distribution.**

AMs and Kernel Machines

- Both **nonparametric**
 - Kernel machines
 - **$O(KD)$** storage and **$O(KD)$** time per inference
 - Associative Memory networks
 - **$O(KD)$** storage and **$O(KDT)$** time per inference
- for a **T** -layer AM network

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

Kernel sum

Opportunities

Can we draw inspiration from the rich literature on kernel machines?

- How can we adapt techniques for **efficient kernel machines** to AM?
- How can we utilize the various **domain-specific kernels** with **unique properties** to design **novel energy functions**, and how do these properties translate to AMs?
- What can we do with this **"mode-finding" capability** of AMs?

Are there other such connections?

How can AMs be used for other ML problems?

Opportunities

Inspiration for Kernel Machines

- Efficiency
- Mode-finding Capabilities
- Novel Energy Functions

Opportunities

Inspiration for Kernel Machines

- **Efficiency**
- Mode-finding Capabilities
- Novel Energy Functions

Memory Capacity & Storage

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

$$\mathbf{T} = \sum_{\mu=1}^K \xi^{\mu} (\xi^{\mu})^{\top}$$

$$O(KD) \rightarrow O(D^2)$$

Classical Hopfield energy

Disentangles number of memories from the number of model parameters needed to store them

(sub)Linear capacity

Memory Capacity & Storage

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

$$\kappa(\mathbf{v}, \boldsymbol{\xi}^{\mu}) = \exp \left(-\beta/2 \|\mathbf{v} - \boldsymbol{\xi}^{\mu}\|^2 \right)$$

Log-sum-exp energy

Need all memories to compute the energy

Exponential capacity

Kernel Feature Maps

MrDAM Memory representation DenseAM

$$\begin{aligned} E_{\beta}(\mathbf{v}; \Xi) &= -Q \left[\sum_{\mu=1}^K \kappa(\mathbf{v}, \xi^{\mu}) \right] \approx -Q \left[\sum_{\mu=1}^K \langle \phi(\mathbf{v}), \phi(\xi^{\mu}) \rangle \right] \\ \text{Explicit feature map} \quad \phi : \mathbb{R}^D &\rightarrow \mathbb{R}^Y \\ \kappa(\mathbf{x}, \mathbf{x}') &\approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \end{aligned}$$
$$\begin{aligned} &= -Q \left[\left\langle \phi(\mathbf{v}), \sum_{\mu=1}^K \phi(\xi^{\mu}) \right\rangle \right] \\ &= -Q [\langle \phi(\mathbf{v}), \mathbf{T} \rangle] \end{aligned}$$

DrDAM Distributed representation DenseAM

Distributed Memories

MrDAM energy: explicit memories

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

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^Y$$

$$\kappa(\mathbf{x}, \mathbf{x}') \approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

$$\kappa(\mathbf{v}, \mathbf{\xi}^{\mu}) = \exp \left(-\beta/2 \|\mathbf{v} - \mathbf{\xi}^{\mu}\|^2 \right)$$

DrDAM energy: distributed memories

$$\tilde{E}_{\beta}(\mathbf{v}; \mathbf{T}) = -Q [\langle \phi(\mathbf{v}), \mathbf{T} \rangle]$$

$$\mathbf{T} = \sum_{\mu=1}^K \phi(\mathbf{\xi}^{\mu}) \in \mathbb{R}^Y$$

Infeasible infinite
dimensional feature map

Random Approximate Feature Maps

Bochner's Theorem:

A shift-invariant kernel is positive definite if and only if it is a Fourier transform of a positive measure

Fourier transform of a positive measure

$$\kappa(\mathbf{z}, \mathbf{z}') = \kappa(\mathbf{z} - \mathbf{z}') = \int_{\mathbb{R}^D} p(\boldsymbol{\omega}) \exp(j \langle \boldsymbol{\omega}, \mathbf{z} - \mathbf{z}' \rangle) d\boldsymbol{\omega}$$

Shift-invariant

Positive measure

For a positive definite shift-invariant kernel, there exists a such a positive measure

Random Features for Kernel Machines

$$\begin{aligned}\kappa(\mathbf{z}, \mathbf{z}') &= \kappa(\mathbf{z} - \mathbf{z}') = \int_{\mathbb{R}^D} p(\boldsymbol{\omega}) e^{j\langle \boldsymbol{\omega}, \mathbf{z} - \mathbf{z}' \rangle} d\boldsymbol{\omega} \\&= \mathbb{E}_{\boldsymbol{\omega} \sim p} e^{j\langle \boldsymbol{\omega}, \mathbf{z} \rangle} \overline{e^{j\langle \boldsymbol{\omega}, \mathbf{z}' \rangle}} \\&= \mathbb{E}_{\boldsymbol{\omega} \sim p} \left\langle \begin{bmatrix} \cos \langle \boldsymbol{\omega}, \mathbf{z} \rangle \\ \sin \langle \boldsymbol{\omega}, \mathbf{z} \rangle \end{bmatrix}, \begin{bmatrix} \cos \langle \boldsymbol{\omega}, \mathbf{z}' \rangle \\ \sin \langle \boldsymbol{\omega}, \mathbf{z}' \rangle \end{bmatrix} \right\rangle \\&\approx \frac{1}{Y} \sum_{i=1}^Y \left\langle \begin{bmatrix} \cos \langle \boldsymbol{\omega}^i, \mathbf{z} \rangle \\ \sin \langle \boldsymbol{\omega}^i, \mathbf{z} \rangle \end{bmatrix}, \begin{bmatrix} \cos \langle \boldsymbol{\omega}^i, \mathbf{z}' \rangle \\ \sin \langle \boldsymbol{\omega}^i, \mathbf{z}' \rangle \end{bmatrix} \right\rangle \\&= \langle \Phi(\mathbf{z}), \Phi(\mathbf{z}') \rangle\end{aligned}$$

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{Y}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}^1, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^1, \mathbf{x} \rangle) \\ \cos(\langle \boldsymbol{\omega}^2, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^2, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}^Y, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}^Y, \mathbf{x} \rangle) \end{bmatrix}$$

$\boldsymbol{\omega}^i \sim p$

Random features
if we know the
positive measure

Random Features for Associative Memories

$$E = -\log \left(\sum_{\mu} \exp \left(-\frac{1}{2} \|\xi^{\mu} - \mathbf{x}\|_2^2 \right) \right)$$

Distributed Memories with Random Features

MrDAM energy: explicit memories

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

$$O(KD) \rightarrow O(Y) = O(D/\epsilon^2)$$

Disentangles number of memories from the number of model parameters needed to store them

DrDAM energy: distributed memories

$$\tilde{E}_{\beta}(\mathbf{v}; \mathbf{T}) = -Q [\langle \Phi(\mathbf{v}), \mathbf{T} \rangle]$$

$$\mathbf{T} = \sum_{\mu=1}^K \Phi(\xi^{\mu}) \in \mathbb{R}^Y$$

$$Y \sim O(D/\epsilon^2)$$

$$|\kappa(\mathbf{z}, \mathbf{z}') - \langle \Phi(\mathbf{z}), \Phi(\mathbf{z}') \rangle| \leq \epsilon$$

Approximation in Energy Descent

MrDAM energy: exact dynamics

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

$$\mathbf{v}^{(t)} \leftarrow \mathbf{v}^{(t-1)} - \eta \nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}^{(t-1)}; \Xi)$$

What is the approximation
in the DenseAM output?

DrDAM energy: approx dynamics

$$\tilde{E}_{\beta}(\mathbf{v}; \mathbf{T}) = -Q [\langle \Phi(\mathbf{v}), \mathbf{T} \rangle]$$

$$\mathbf{T} = \sum_{\mu=1}^K \Phi(\xi^{\mu}) \in \mathbb{R}^Y$$

$$\tilde{\mathbf{v}}^{(t)} \leftarrow \tilde{\mathbf{v}}^{(t-1)} - \eta \nabla_{\mathbf{v}} \tilde{E}_{\beta}(\mathbf{v}^{(t-1)}; \mathbf{T})$$

$$\left\| \mathbf{v}^{(T)} - \tilde{\mathbf{v}}^{(T)} \right\| \leq ?$$

Approximation in Energy Descent

MrDAM energy: exact dynamics

$$\mathbf{v}^{(t)} \leftarrow \mathbf{v}^{(t-1)} - \eta \nabla_{\mathbf{v}} E_{\beta}(\mathbf{v}^{(t-1)}; \Xi)$$

Random feature approx bound

$$|\kappa(\mathbf{z}, \mathbf{z}') - \langle \Phi(\mathbf{z}), \Phi(\mathbf{z}') \rangle| \leq C_1 \sqrt{D/Y}$$

DrDAM energy: approx dynamics

$$\tilde{\mathbf{v}}^{(t)} \leftarrow \tilde{\mathbf{v}}^{(t-1)} - \eta \nabla_{\mathbf{v}} \tilde{E}_{\beta}(\mathbf{v}^{(t-1)}; \mathbf{T})$$

Sufficiently small step-size

$$\eta \leq \frac{C_2}{T(1 + 2K\beta \exp(\beta/2))}$$

Initial energy

$$\left\| \mathbf{v}^{(T)} - \tilde{\mathbf{v}}^{(T)} \right\| \leq \frac{C_1 C_2 \exp(\beta(E_{\beta}(\mathbf{v}; \Xi) - 1/2))}{\beta(1 - C_2)}$$

Dense Associative Memory Through the Lens of Random Features

Benjamin Hoover
IBM Research; Georgia Tech
benjamin.hoover@ibm.com

Duen Horng Chau
Georgia Tech
polo@gatech.edu

Hendrik Strobelt
IBM Research; MIT-IBM
hendrik.strobelt@ibm.com

Parikshit Ram
IBM Research
parikshit.ram@ibm.com

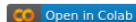
Dmitry Krotov
IBM Research
krotov@ibm.com



Associative Memory Tutorial
Getting started
lib
Pokemon Sprites
tutorial
Binary Dense Storage
Energy Transformer
Memory and Diffusion
Distributed Memory

Distributed Memory

Random Features enable Dense Associative Memory to store patterns in a distributed manner across a large number of neurons.

 Open in Colab

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_{Ξ} 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 $\Xi = \{\xi^{\mu}, \mu \in [K]\}$, thus giving us a model $f_{\mathbf{T}}$ of size $O(Y)$.

On this page

[Exact Energy Function](#)

Visualizing the Energy in 2D

Minimizing the Energy via
Gradient Descent

Viewing Energy as a Kernel Sum

Approximating the Energy with
Random Features

 Report an issue

Other Formats

 CommonMark



Opportunities

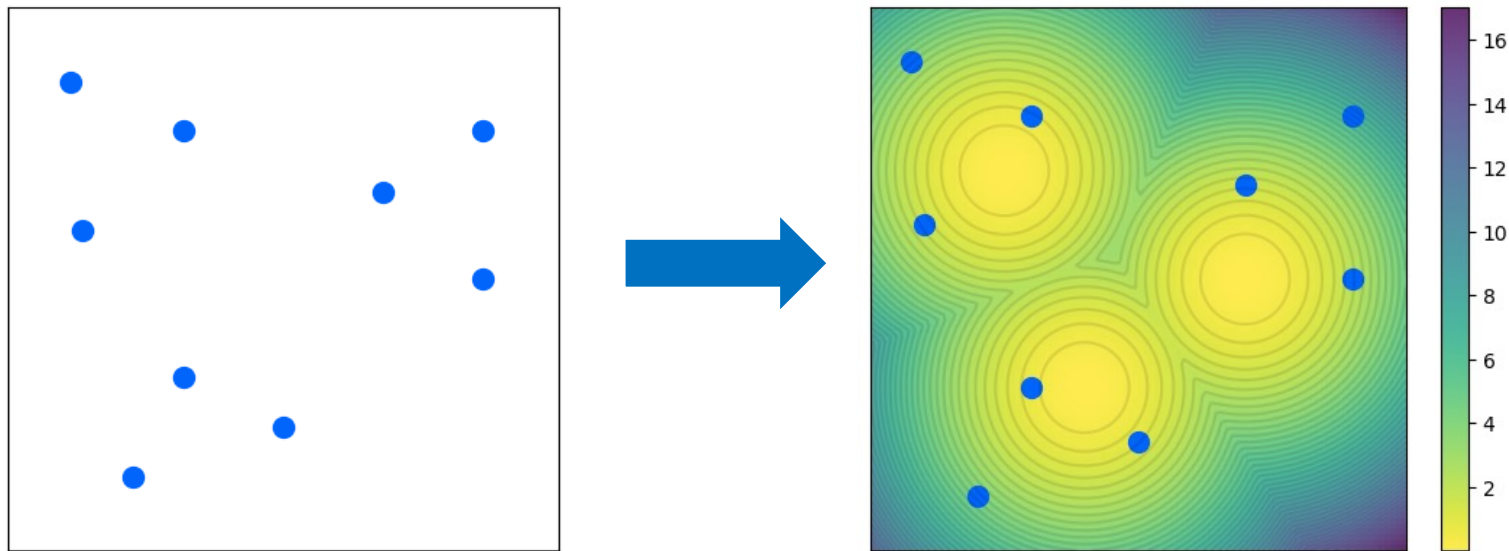
Inspiration for Kernel Machines

- Efficiency
- **Mode-finding Capabilities**
- Novel Energy Functions

Contractive Layer

- Layer operates on inputs independently
- But can exhibit collective contraction

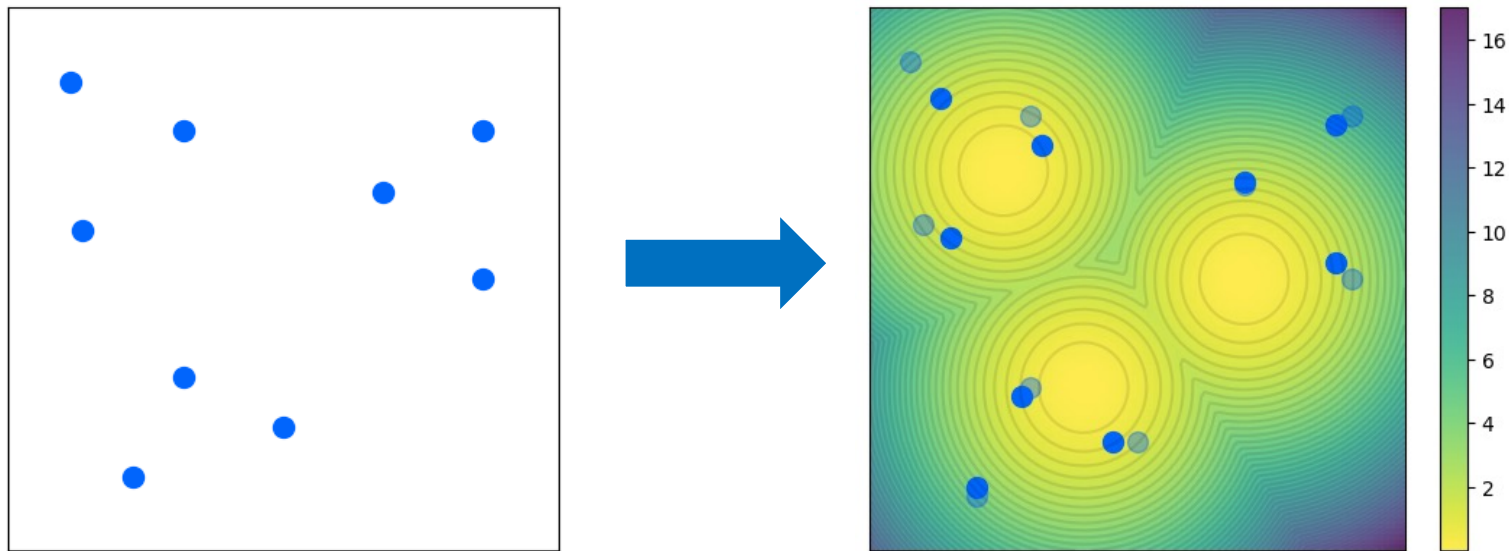
$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$



Contractive Layer

- Layer operates on inputs independently
- But can exhibit collective contraction

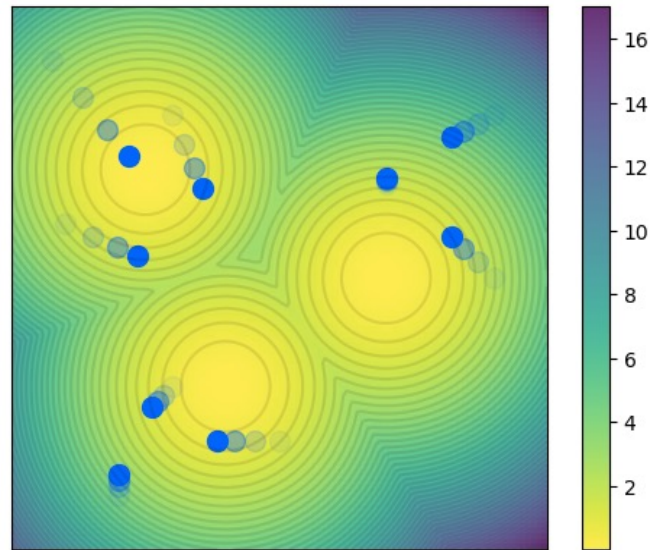
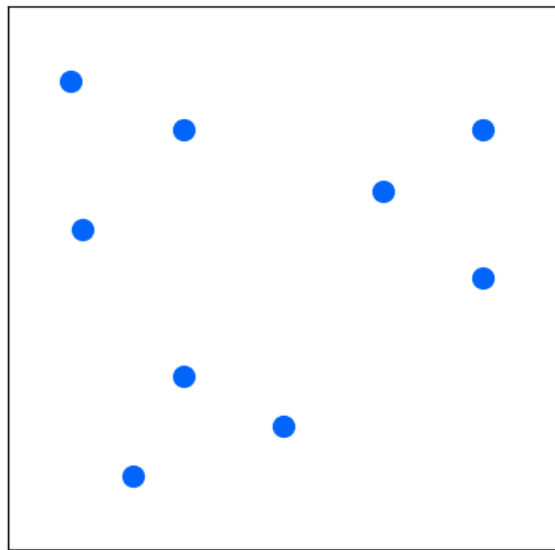
$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$



Contractive Layer

- Layer operates on inputs independently
- But can exhibit collective contraction by **contracting towards modes**

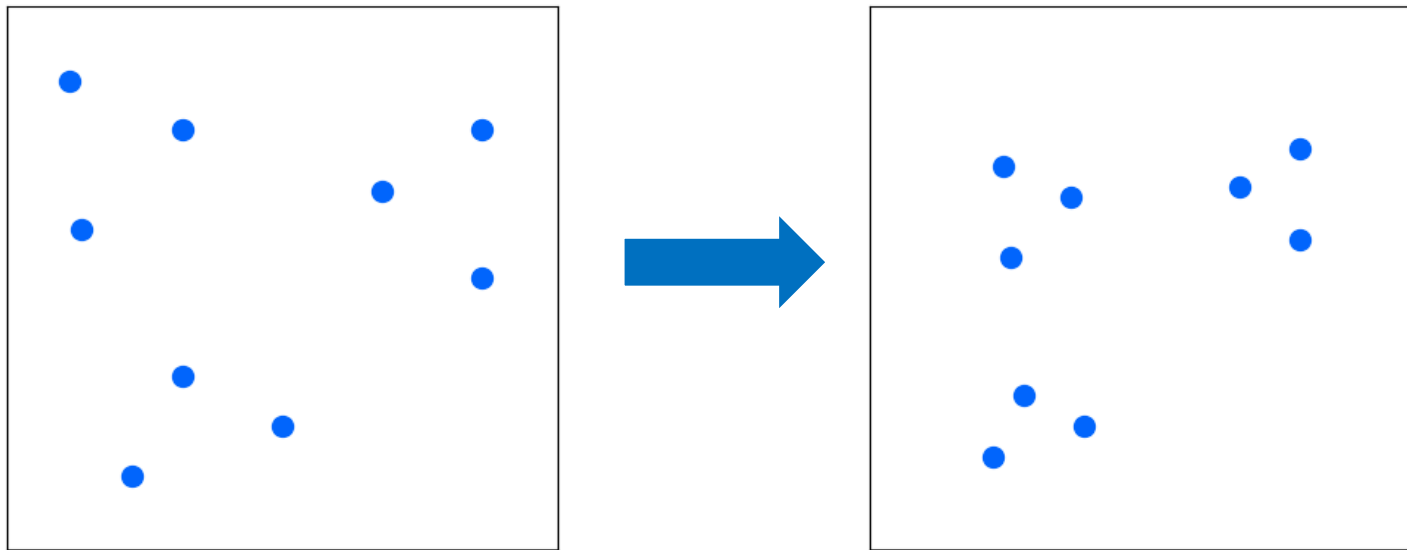
$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$



Contractive Layer

- Layer operates on inputs independently
- But can exhibit collective contraction by **contracting towards modes**

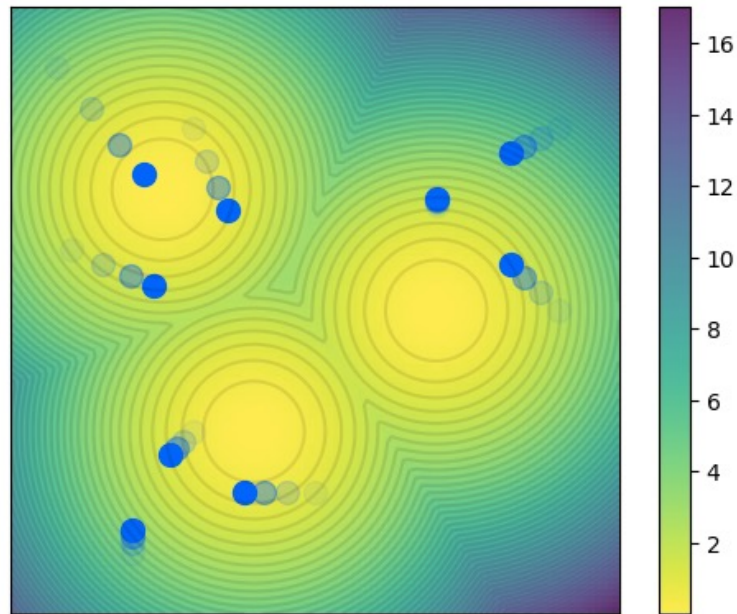
$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$



Contractive Layer

- Parameters of the AM control locations of the minima (modes)
- Thus, where the inputs contract towards

$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$



Clustering

Discrete k -means clustering

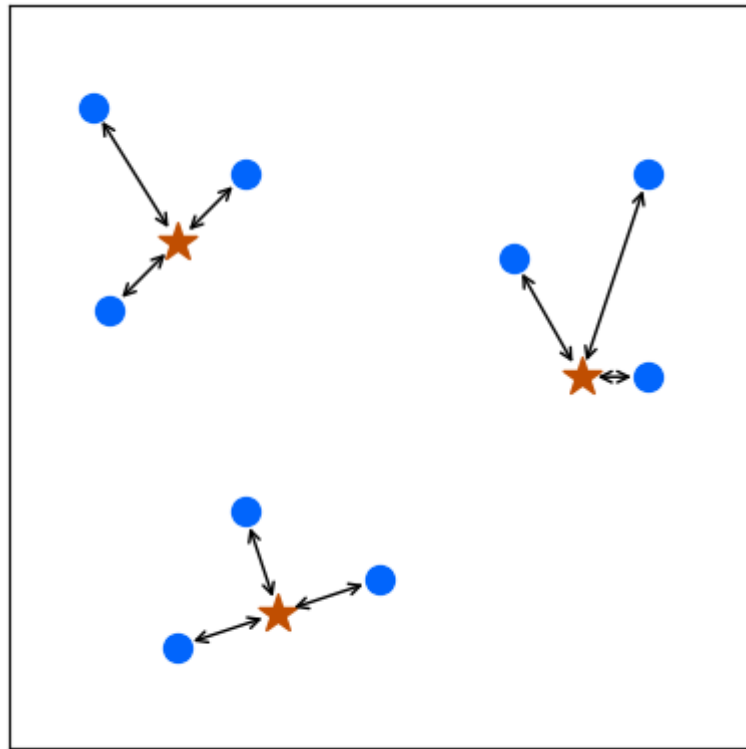
Points to cluster

$$\min_{\mathbf{c}^1, \dots, \mathbf{c}^k \in \mathbb{R}^D} \sum_{i=1}^m \min_{j \in [k]} \|\mathbf{x}^i - \mathbf{c}^j\|^2$$

Centers

Distance to
closest center

**Discrete optimization – we need
discrete assignments to clusters**



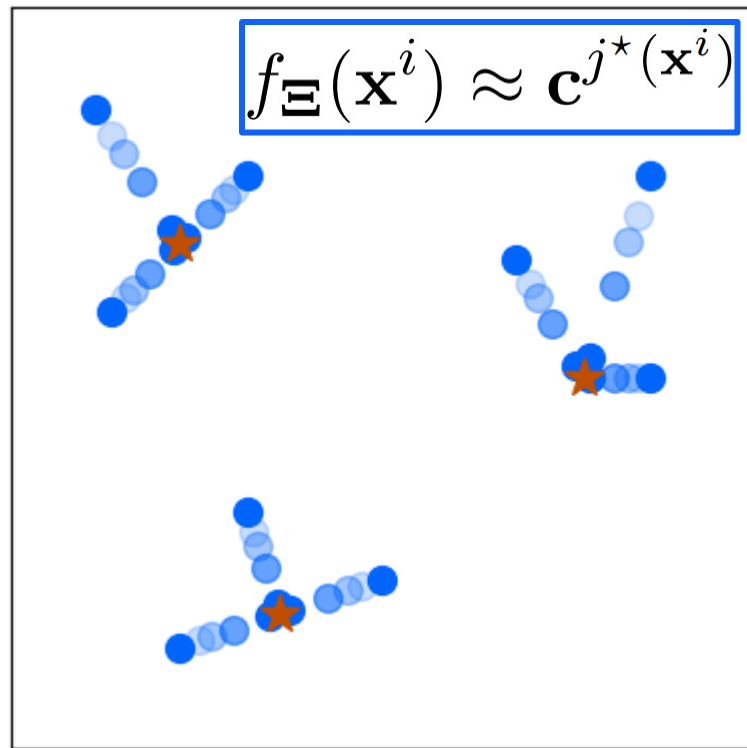
Clustering with Associative Memories

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

Main idea: Use contraction to emulate discrete assignment

$$\min_{\Xi} \sum_{i=1}^m \|\mathbf{x}^i - f_{\Xi}(\mathbf{x}^i)\|^2$$

Differentiable discrete clustering objective



End-to-end Differentiable Clustering with Associative Memories

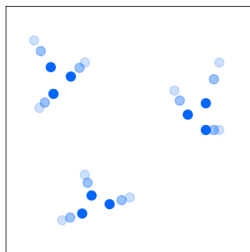
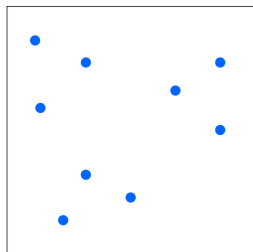
Bishwajit Saha¹ Dmitry Krotov² Mohammed J. Zaki¹ Parikshit Ram^{2,3}



Deep Clustering with Associative Memories

Main idea: Use contraction to learn a clustered latent space

$$\mathbf{x} \xrightarrow{\text{encode}} e_{\varphi}(\mathbf{x}) \xrightarrow{\text{contract}} f_{\Xi} \circ e_{\varphi}(\mathbf{x}) \xrightarrow{\text{decode}} d_{\vartheta} \circ f_{\Xi} \circ e_{\varphi}(\mathbf{x})$$



$$\min_{\varphi, \Xi, \vartheta} \sum_{\mathbf{x}} \|\mathbf{x} - d_{\vartheta} \circ f_{\Xi} \circ e_{\varphi}(\mathbf{x})\|^2$$

Single differentiable objective handling both **fidelity of learned representations** and the **collective clustered structure**

← Go to **ICLR 2025 Workshop NFAM** homepage

Deep Clustering with Associative Memories

Bishwajit Saha, Dmitry Krotov, Mohammed J Zaki, Parikshit Ram



Opportunities

Inspiration for Kernel Machines

- Efficiency
- Mode-finding Capabilities
- **Novel Energy Functions**

Energy Functions

$$f_{\Xi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$$

$$\Xi = [\xi^1, \xi^2, \dots, \xi^K], \xi^\mu \in \mathbb{R}^D$$

Model is parameterized
with **stored patterns**

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

Energy function
defined by the kernel

- Gaussian kernel -- the **log-sum-exp or LSE energy**
- Kernel uses exponential separation with exponential capacity
- Is that enough to make it a "good" kernel function?

Insights from Density Estimation

$$\Xi = [\xi^1, \xi^2, \dots, \xi^K], \xi^\mu \in \mathbb{R}^D \quad \xi^\mu \sim p_{\text{data}}, \mu \in \llbracket K \rrbracket$$

Kernel density estimate or KDE given samples from a distribution

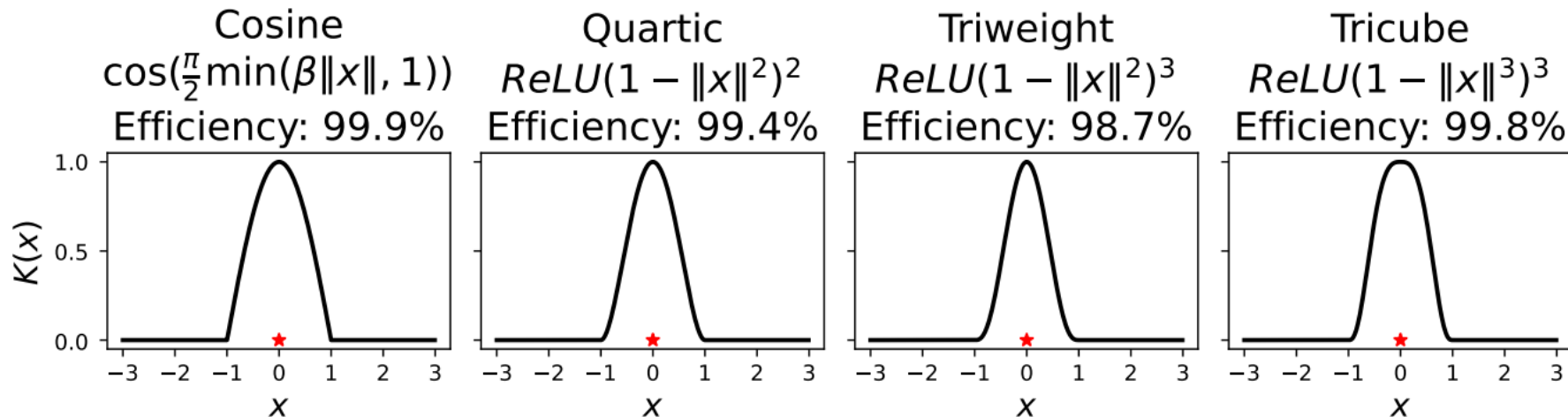
$$\hat{p}_h(\mathbf{v}; \Xi) = \frac{1}{Kh} \sum_{\mu=1}^K \kappa \left(\frac{\mathbf{v} - \xi^\mu}{h} \right)$$

Mean Integrated Squared Error

$$\text{MISE}(h) = \mathbb{E} \left[\int_v (\hat{p}_h(v; \Xi) - p_{\text{data}}(v))^2 dv \right] \sim O \left(\left(\underbrace{\sqrt{\int_z z^2 \kappa(z) dz}}_{\text{Bias}} \underbrace{\int_z \kappa(z)^2 dz}_{\text{Variance}} \right)^{\frac{4}{5}} \right)$$

Kernels and their Efficiencies

$$\sqrt{\int_{\mathbf{z}} \mathbf{z}^2 \kappa(\mathbf{z}) d\mathbf{z} \int_{\mathbf{z}} \kappa(\mathbf{z})^2 d\mathbf{z}}$$



For density estimation

- Gaussian kernel is not the best; many other better
- Epanechnikov kernel known to be optimal

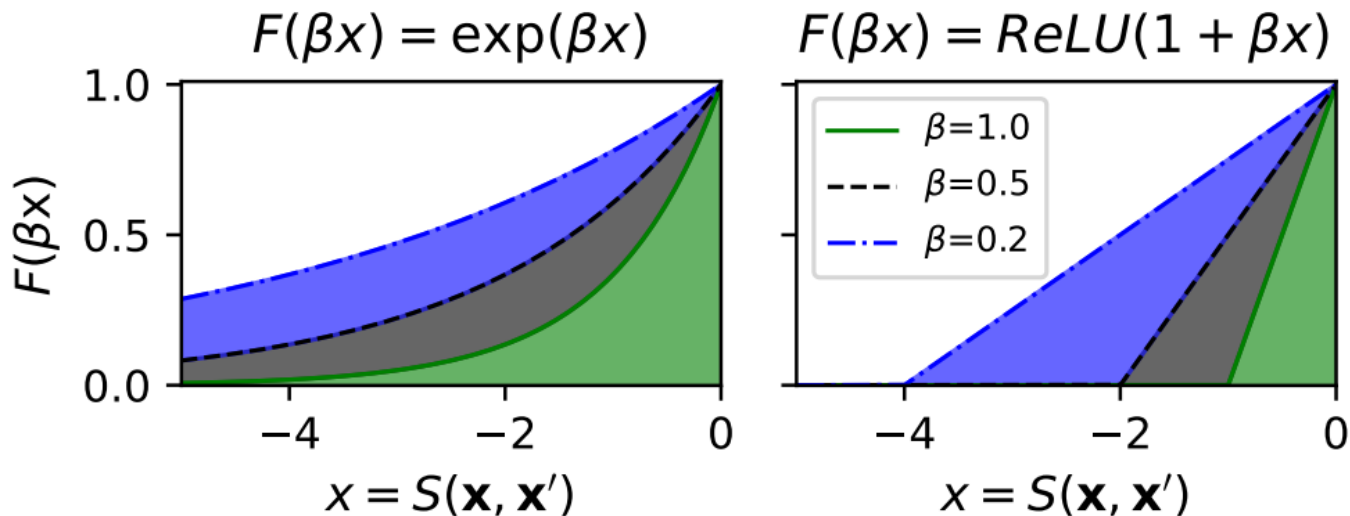
Log-Sum-Shifted-ReLU or LSR Energy

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

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

- Exponential capacity *without* exponential separation function
- Simultaneously retrieves memories **and** generates many new local minima

Log-Sum-Shifted-ReLU or LSR Energy

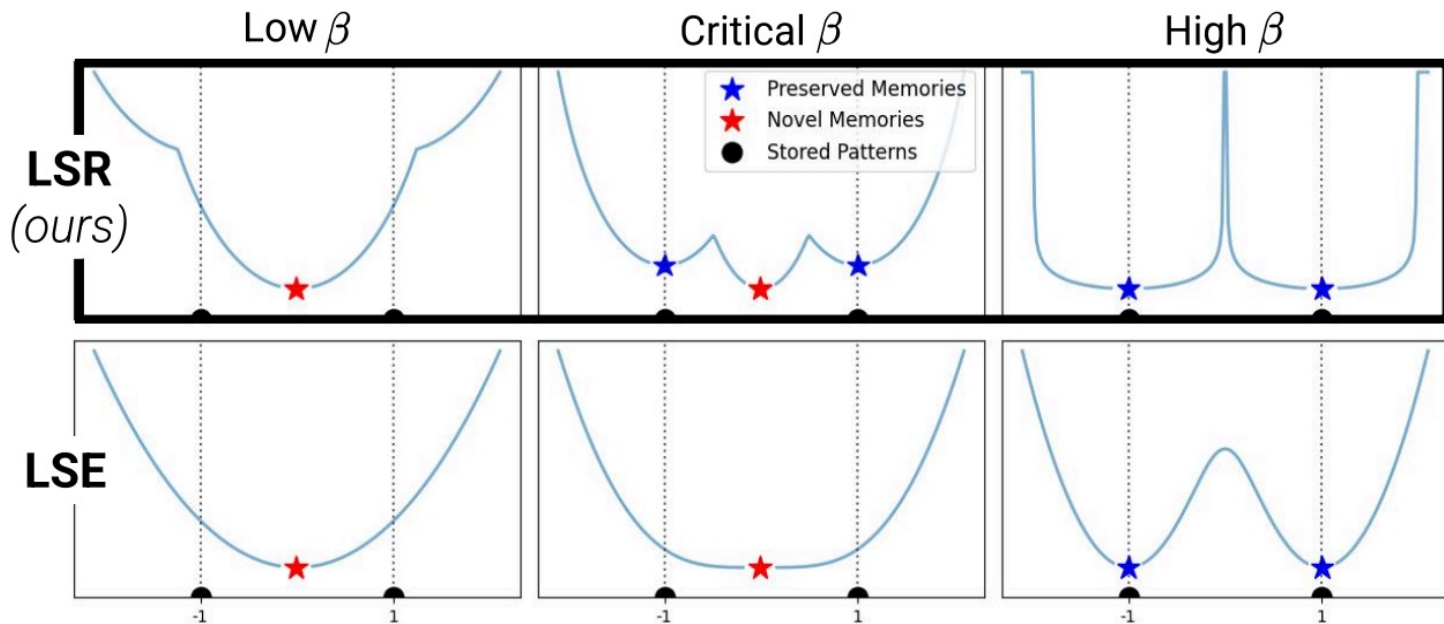


- Exponential capacity *without* exponential separation function
- Simultaneously retrieves memories **and** generates many new local minima

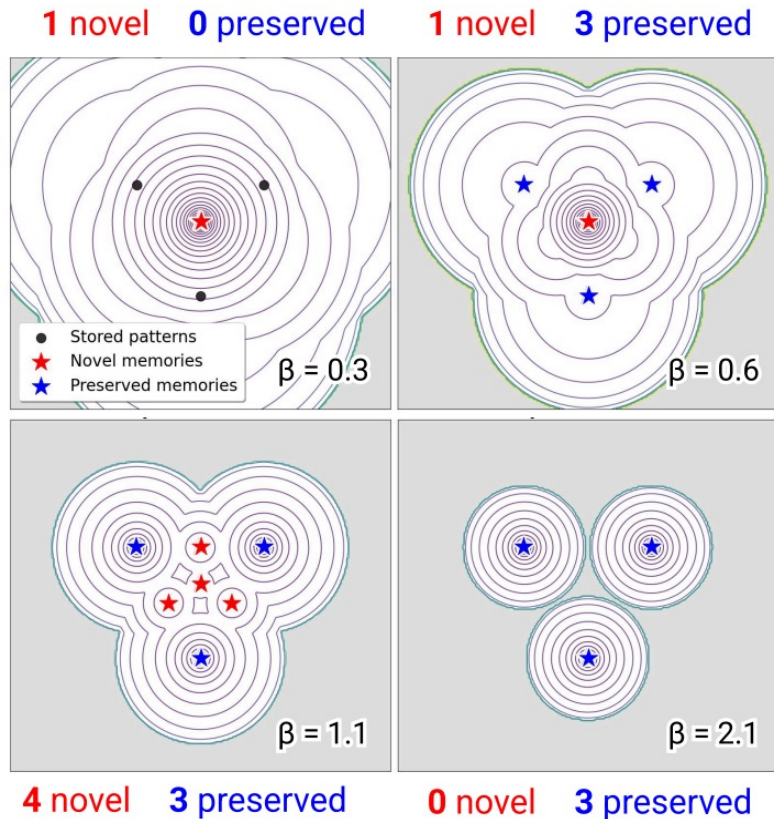
Epanechnikov Energy and Emergent Memories

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

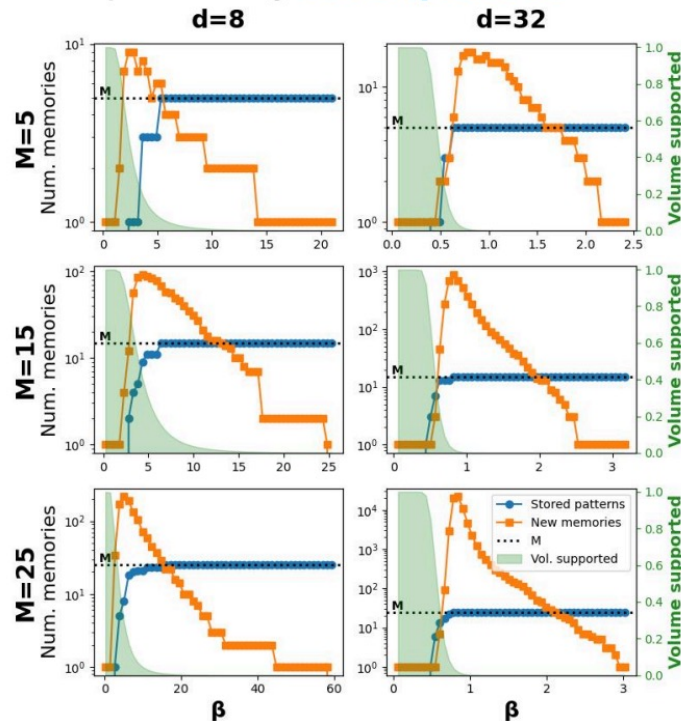
LSE can do only one or the other.



Epanechnikov Energy and Emergent Memories





LSR Energy creates **novel memories** while preserving **stored patterns**



Epanechnikov Energy and Emergent Memories

22 / 24 **preserved** memories

9 0 4 8 6 2 6 8
8 0 0 9 3 3 2 5
9 9 1 5  8  0

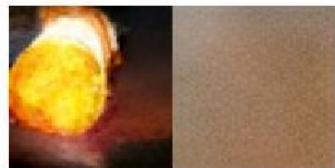
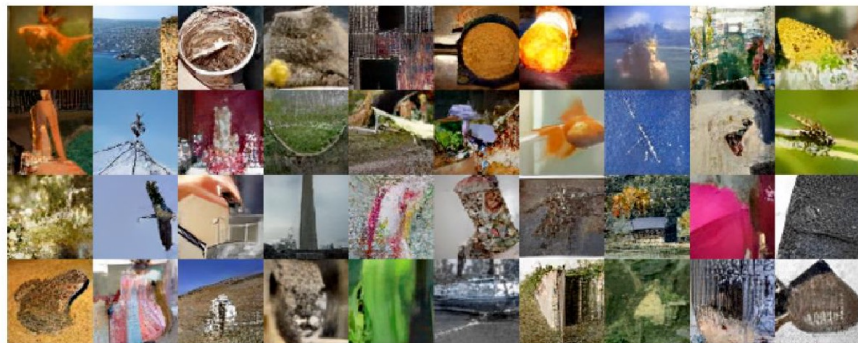
0 2 6 1 0 2

46 **emergent** memories

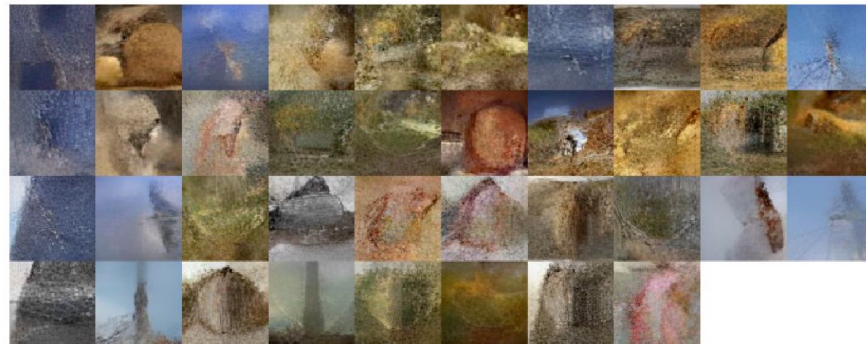
1 8 1 5 6 8 1 0 7 3 1 3
2 8 2 6 6 1 1 3 8 9 9 0
8 3 9 9 6 0 8 3 3 9 8 3
8 8 8 8 2 3 8 9 8 8

Epanechnikov Energy and Emergent Memories

40 / 40 preserved memories



38 emergent memories



[← Go to ICLR 2025 Workshop NFAM homepage](#)

Dense Associative Memory with Epanechnikov energy

Benjamin Hoover, Krishna Balasubramanian, Dmitry Krotov, Parikshit Ram



**How to think about
Associative Memory
Networks as a
Machine Learning
Model?**



**How to use them for
standard Machine
Learning tasks?**

Chapter 5

Associative Memory:
A Machine Learning Model