

1. Introduction

Background- Modern Hopfield Networks (MHNs)

MHNs are energy-based models that can retrieve a semantically relevant data for a given query data.

Issue of MHN

For an out-of-distribution (OOD) query, MHN inevitably associates an inappropriate in-distribution (ID) data.

Contributions

- To address this, we propose **Rectified Lagrangian (ReLag)** that explicitly incorporates OOD queries with a specially designed attractor in the dynamical system of MHNs.
- ReLag-based MHNs yield higher OOD detection performance over existing energy-based methods on average using nine image datasets.

2. Large-Memory MHNs [D. Krotov+, ICLR2021; M. Widrich+, NeurIPS2021]

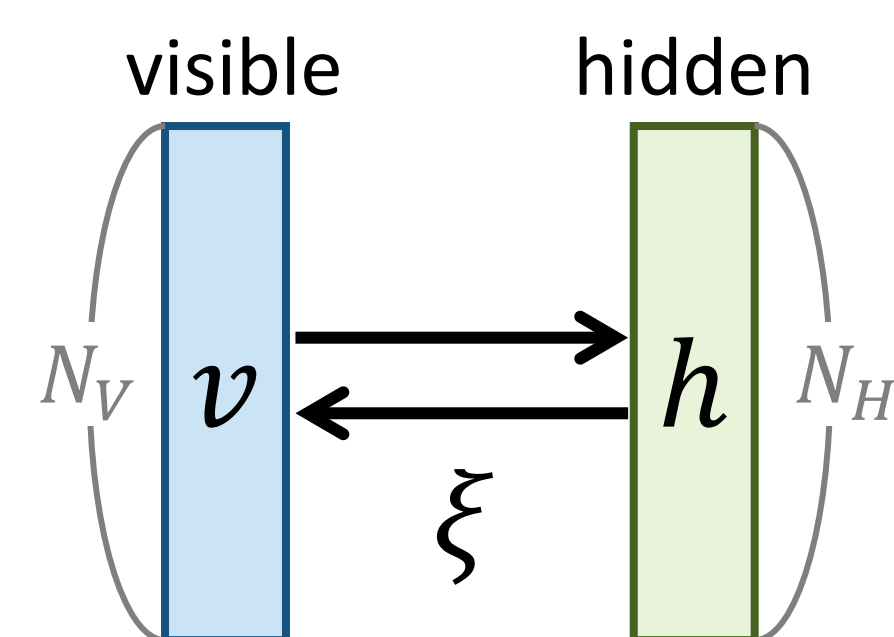
General Formulation of Energy Function

$$E_{\xi}(v, h) = v^{\top} \frac{\partial L_V(v)}{\partial v} - L_V(v) + h^{\top} \frac{\partial L_H(h)}{\partial h} - L_H(h) - \frac{\partial L_H(h)}{\partial h}^{\top} \xi \frac{\partial L_V(v)}{\partial v}$$

Energy associated with visible layer

Energy associated with hidden layer

Interaction energy between visible & hidden layers



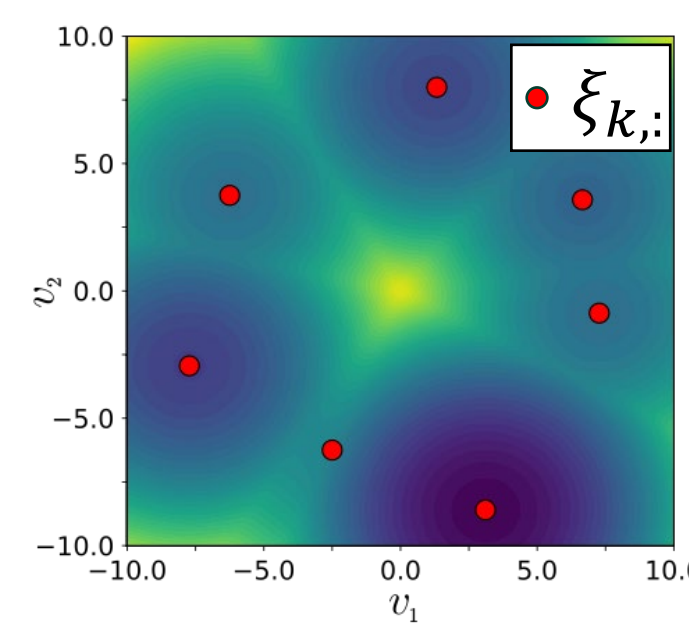
Large-Memory Lagrangian Functions

$$L_V = \frac{1}{2} \sum_{i=1}^{N_V} v_i^2 \quad L_H = \frac{1}{\beta} \log \sum_{k=1}^{N_H} e^{\beta h_k} \quad (\beta > 0)$$

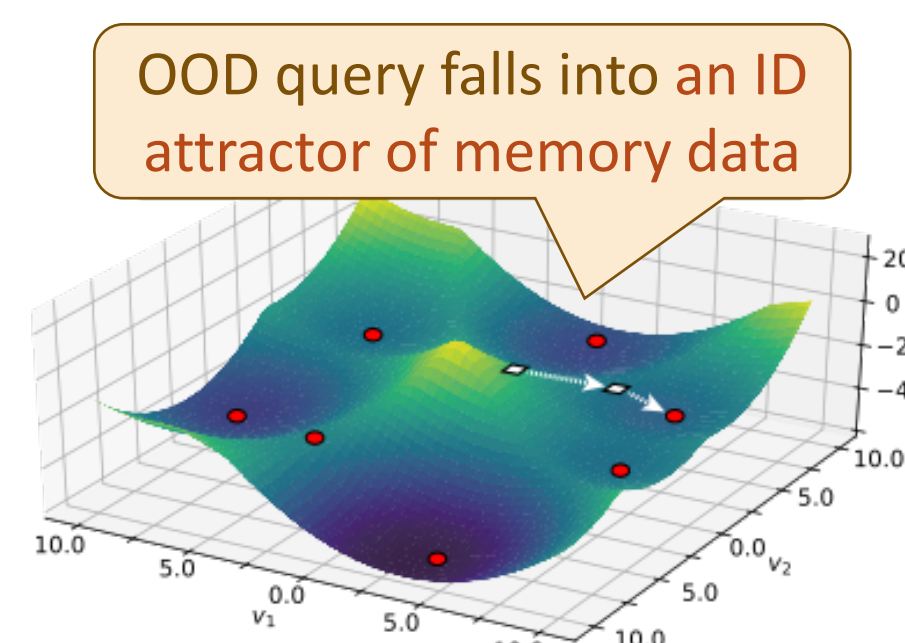
Proven that **log-sum-exp** assures exponentially-ordered memory capacity.

Discrete Update Rule

$v \leftarrow \xi^{\top} \text{softmax}(\beta \xi v)$ Converges to an **attractor**: $v \rightarrow \xi_{k,:}$ (under $0 < \beta \ll 1$)



Energy landscape with attractors



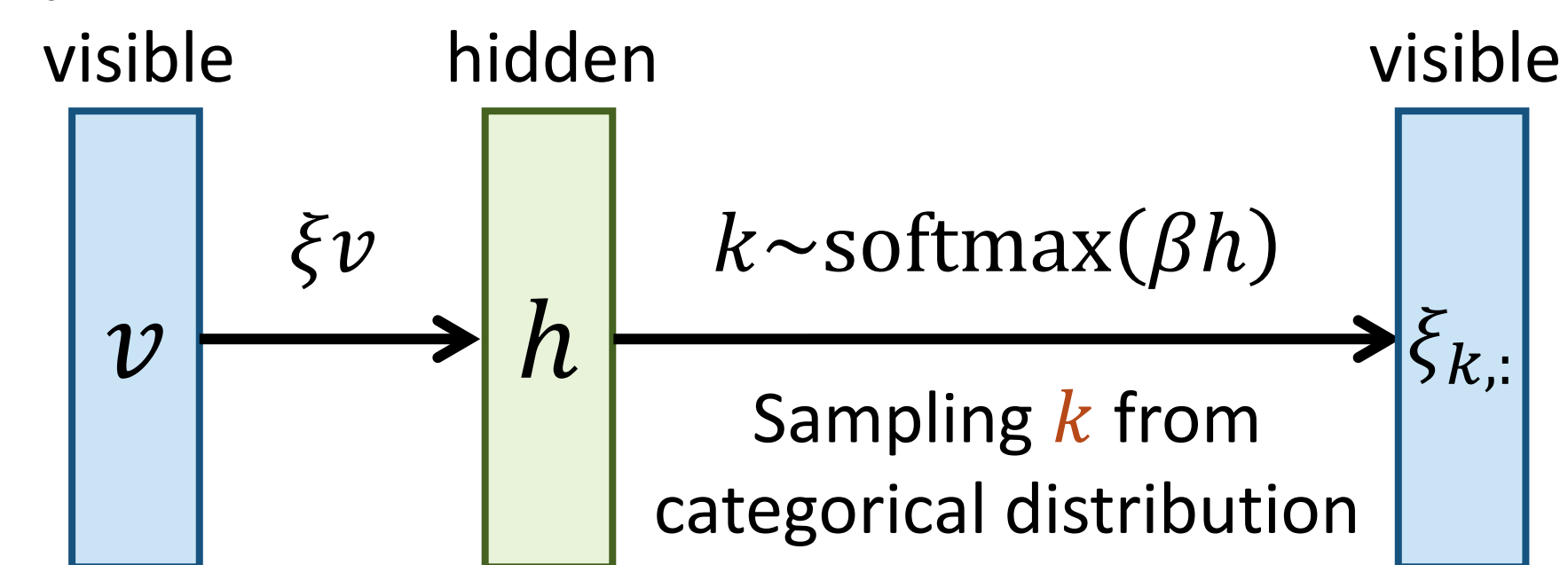
Issue: An OOD query is associated with an irrelevant attractor.

3. Proposed Method: ReLag

Probability-Density Aware MHNs

We optimize ξ to express joint prob. of v & the row index k as: $p(v, k) \propto e^{\beta \xi_{k,:}^{\top} v} \propto \text{softmax}(\beta \xi^{\top} v)_k$

Step 1) Unrolling RNN



The attractor dynamics is virtually acquired when the input is reconstructed.

Step 2) Optimizing probabilistic AE $\xi^* = \arg \max_{\xi} \prod_{v_j \in \mathcal{D}} \prod_{\tilde{v} \sim p(v)} \mathcal{N}(\tilde{v} | v_j; I_{M \times M})$

The probability density for v is given by: $\log p(v) = \log \sum_{k=1}^{N_H} e^{\beta \xi_{k,:}^{\top} v} + \text{const}$

hidden-layer Lagrangian L_H

Now, **log-sum-exp** is connected to $p(v)$.

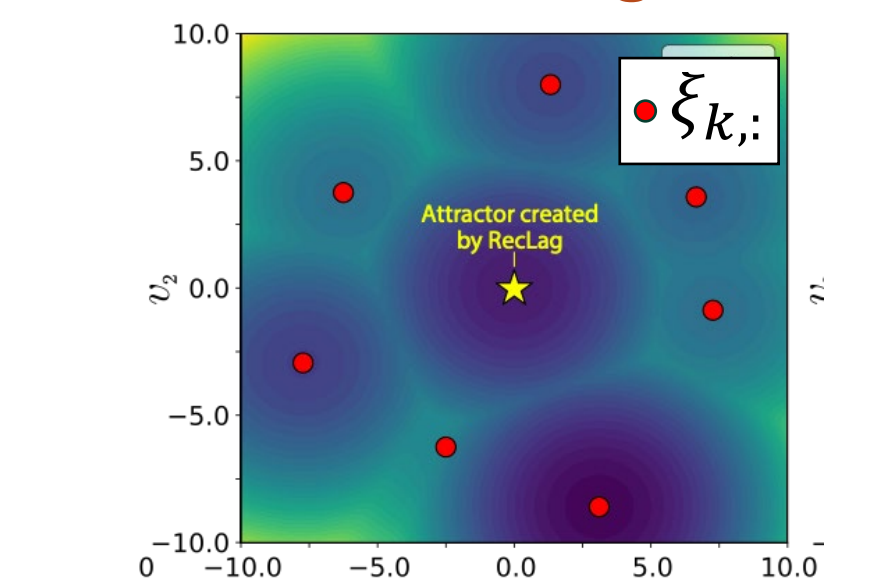
Rectified Lagrangian Containing a Trivial Attractor

We replace the hidden-layer Lagrangian with:

$$L_H^{\text{Rec}}(h) = \max \left[\frac{1}{\beta} \log \gamma, \frac{1}{\beta} \log \sum_{k=1}^{N_H} e^{\beta h_k} \right] \quad \gamma > 0$$

rectification

log-sum-exp

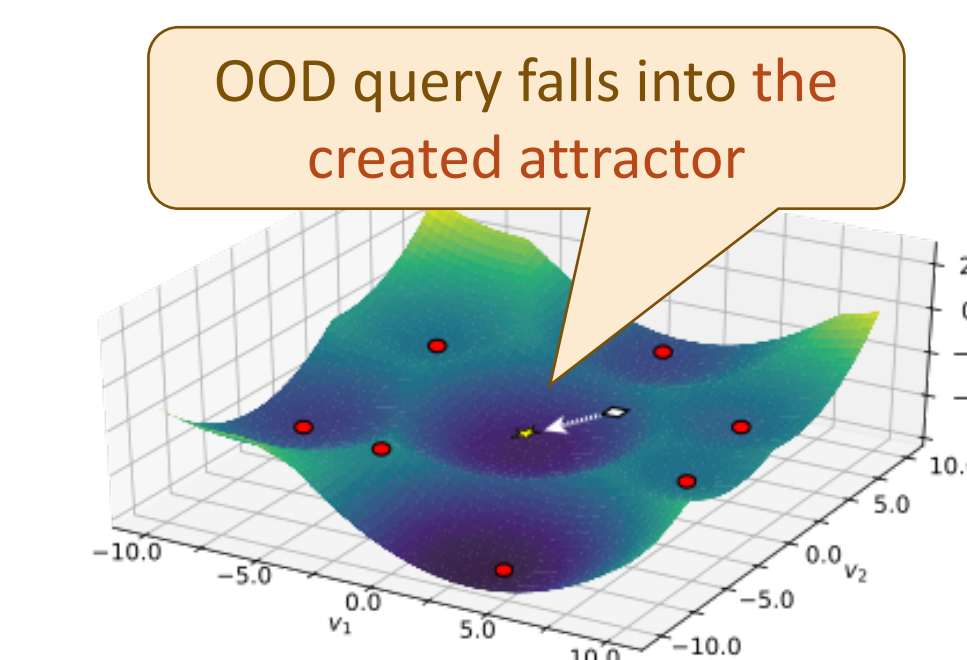


Energy landscape with attractors

Modified Update Rule

$$v \leftarrow \begin{cases} \xi^{\top} \text{softmax}(\beta \xi v), & L_H(h) \geq \beta^{-1} \log \gamma \\ 0, & L_H(h) < \beta^{-1} \log \gamma \end{cases}$$

attractor absorbing OOD queries



OOD query falls into the created attractor

An OOD query is associated with the created attractor.

4. Evaluation

Outline

OOD detection performance is evaluated against 5 baselines using 9 datasets.

Result

- ReLag-based MHNs yield higher OOD detection performance over baselines on average.
- ReLag-based MHNs also showed the highest accuracy in 23 out of 27 individual experimental settings.

[D. Hendrycks+, ICLR2017; W. Liu+, NeurIPS2020; Y. Sun+, NeurIPS2021; J Zhang+, ICLR2023]

Evaluation

Table: OOD detection performance as FPR95(%) ↓

Method	SVHN	LSUN-C	LSUN-R	iSUN	Places	DTD	TIN	SUN	iNaturalist	Average
ResNet18	MSP	76.34	27.52	36.54	34.84	20.55	30.65	45.82	22.89	12.62
	Energy	56.05	8.10	11.60	9.10	3.18	16.98	25.47	3.27	3.47
	ReAct	59.47	7.57	12.52	10.13	2.93	16.86	27.61	3.27	3.80
	MHE	17.59	9.20	7.68	4.74	0.33	8.96	15.86	0.00	2.35
	SHE	17.45	9.22	7.69	4.77	0.33	8.99	15.84	0.00	2.38
ReLag	18.12	6.40	4.60	2.67	0.28	6.82	12.09	0.00	1.68	5.85
	± 2.02	± 0.25	± 0.12	± 0.47	± 0.02	± 0.13	± 0.25	± 0.00	± 0.04	± 0.24
ResNet34	MSP	59.86	28.26	32.06	31.69	33.61	43.28	45.56	32.43	32.95
	Energy	30.51	6.84	9.43	8.47	9.32	23.74	25.16	8.99	10.86
	ReAct	45.86	14.37	14.09	13.28	15.83	29.73	31.60	15.53	11.98
	MHE	6.20	6.17	4.40	2.94	2.34	14.32	15.86	0.54	4.91
	SHE	6.14	6.20	4.45	3.01	2.36	14.32	15.93	0.54	4.92
ReLag	5.19	5.60	2.85	2.11	2.31	12.04	11.71	0.33	4.14	5.14
	± 0.24	± 0.07	± 0.05	± 0.05	± 0.03	± 0.07	± 0.23	± 0.11	± 0.08	± 0.08
WRN40-2	MSP	41.52	44.43	38.47	39.70	33.84	35.80	51.52	34.88	27.69
	Energy	15.35	17.77	14.98	17.45	10.58	19.71	36.75	9.54	8.95
	ReAct	18.83	19.93	18.25	20.68	11.98	21.67	42.02	11.44	13.26
	MHE	5.40	14.60	12.03	11.48	2.90	10.99	27.28	0.82	1.83
	SHE	5.25	14.39	13.18	12.39	2.83	10.98	28.35	0.82	1.84
ReLag	5.75	7.37	8.44	8.01	2.63	9.75	22.62	1.06	1.67	7.47
	± 0.12	± 0.18	± 0.17	± 0.15	± 0.05	± 0.10	± 0.34	± 0.09	± 0.05	± 0.85

Rectified Lagrangian for Out-of-Distribution Detection in Modern Hopfield Networks

Ryo Moriai¹, Nakamasa Inoue¹, Masayuki Tanaka¹, Rei Kawakami¹, Satoshi Ikehata^{1,3}, Ikuro Sato^{1,2}¹Tokyo Institute of Technology, Japan, ²Denso IT Laboratory, Japan, ³National Institute of Informatics, Japan

Introduction

Modern Hopfield Networks (MHNs)

Energy-based models that can retrieve a relevant data for a given query data.

For out-of-distribution (OOD) query, MHN does not retrieve in-distribution (ID) data, which is totally different from the query.

To address this, we propose Rectified Lagrangian based MHNs. We introduce a new Lagrangian that explicitly penalizes OOD queries with a specially designed Rectified Lagrangian in the dynamical system of MHNs. Our Rectified Lagrangian-based MHNs demonstrated higher OOD detection performance over existing energy-based methods on average using nine datasets.

2. Modern Hopfield Network (MHN) [D. Krotov+ 2021, Widrich et al. 2021]

Formulation

 $v \in \mathbb{R}^{N_V}$: visible layer neuron $h \in \mathbb{R}^{N_H}$: hidden layer neuron,

 β : real hyperparameter,

 $\Xi \in \mathbb{R}^{N_H \times N_V}$: matrix parameter,

 ξ_μ : μ -th column vector component of a matrix Ξ

Energy

$$E = \sum_{i=1}^{N_V} v_i g_i(v) - L_V(v) + \sum_{\mu=1}^{N_H} h_\mu f_\mu(h) - L_H(h) - f(h)^T \Xi g(v)$$

Activation functions

$$f(h) = \frac{\partial L_H(h)}{\partial h} \quad g(v) = \frac{\partial L_V(v)}{\partial v}$$

Hidden-layer Lagrangian

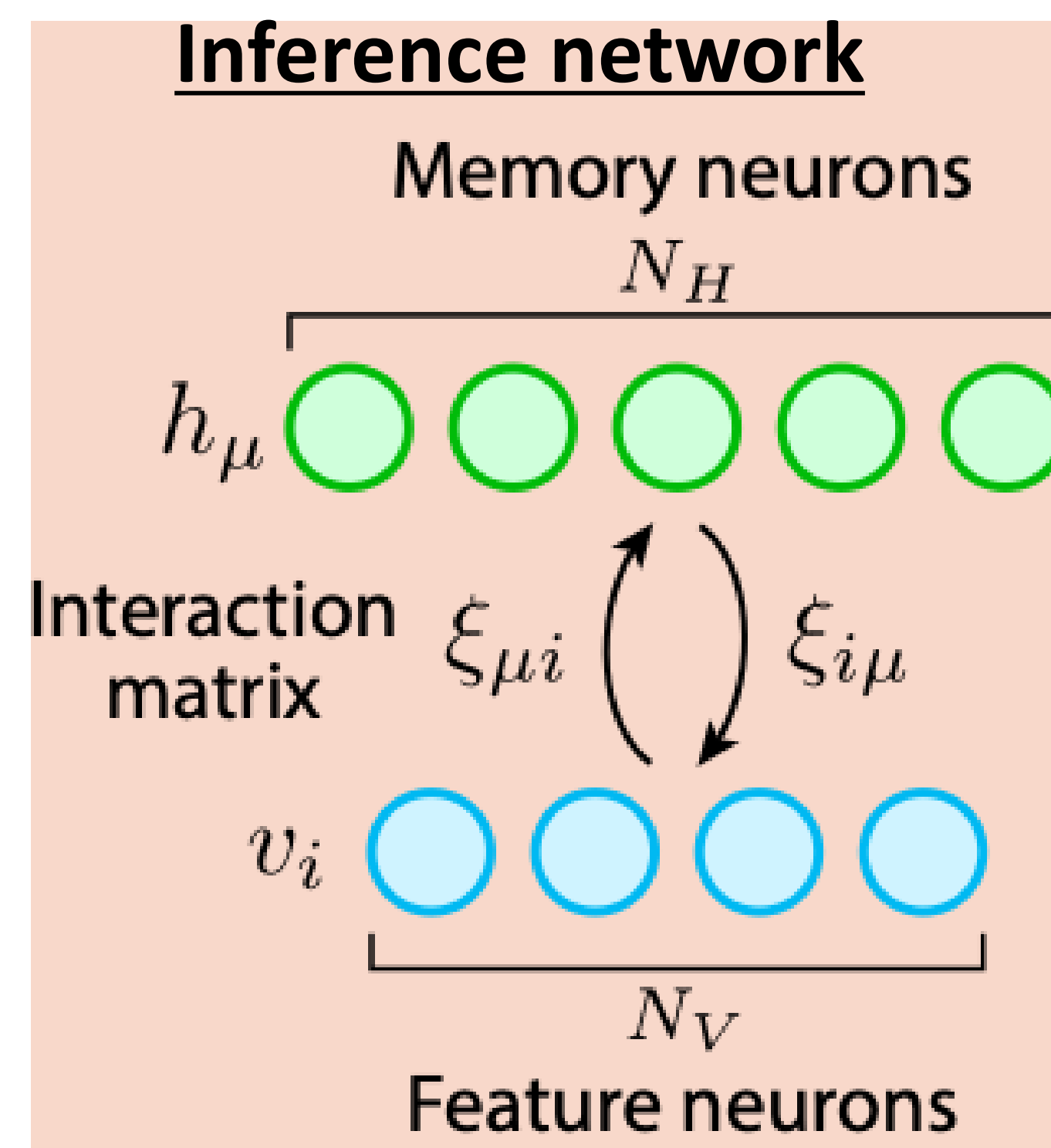
$$L_h = \frac{1}{\beta} \log \left(\sum_{\mu=1}^{N_H} e^{\beta h_\mu} \right)$$

Visible Lagrangian

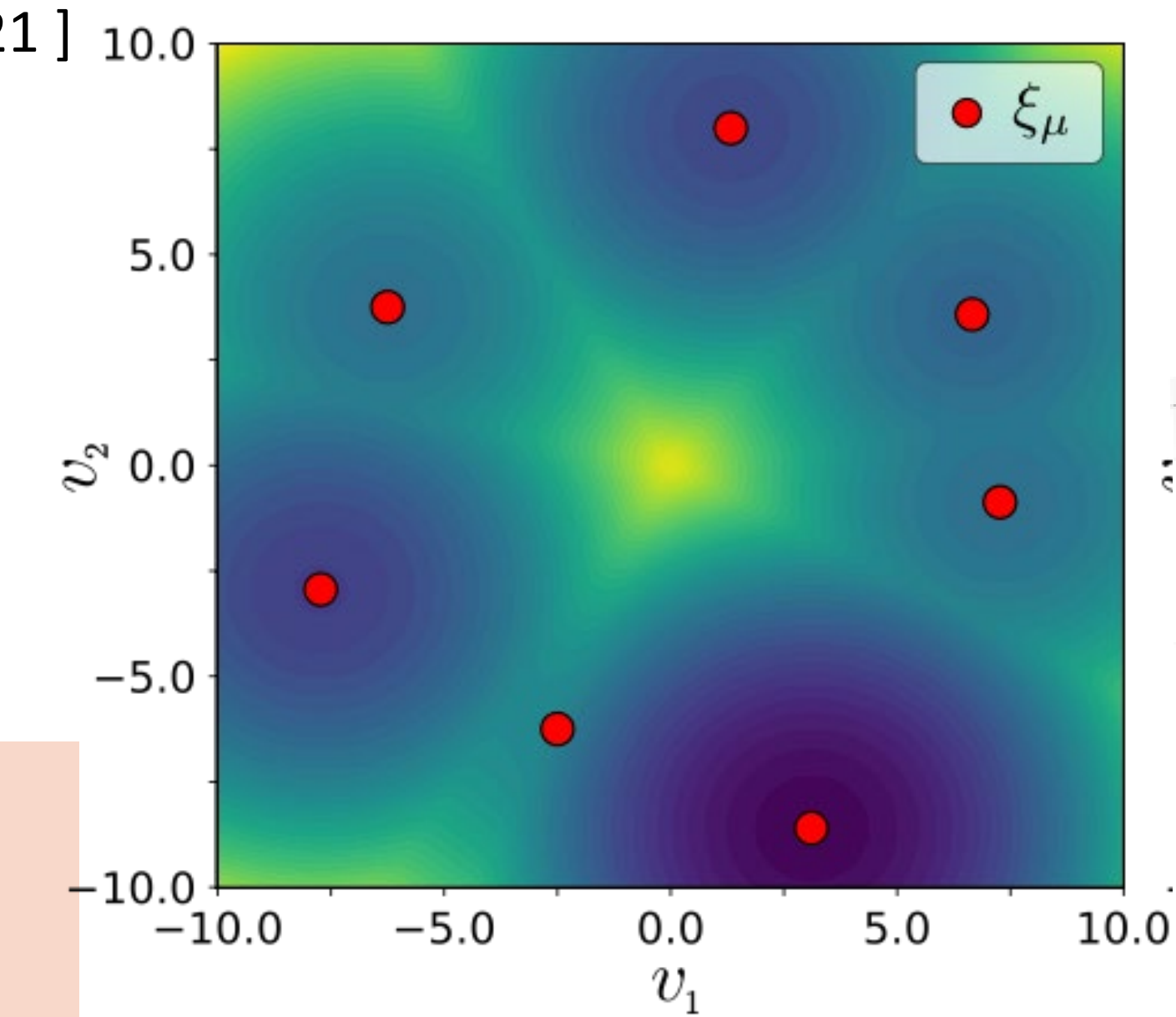
$$L_v = \frac{1}{2} \sum_{i=1}^{N_V} v_i^2$$

Discrete Update rule

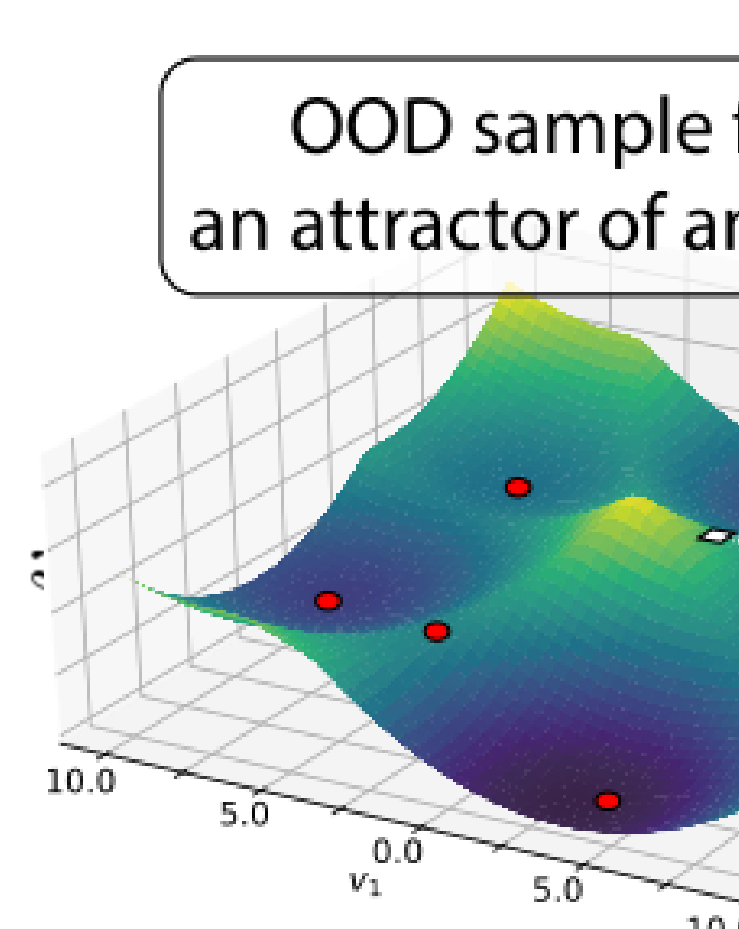
$$v^{k+1} = \Xi^T \text{softmax}(\beta \Xi v^k)$$



Energy Map



OOD sample



4. Experiment

Experiment

Method	SVHN	LSUN-C	LSUN-R	iSUN	Places	DTD	TIN	SUN
ResNet18	MSP	76.34	27.52	36.54	34.84	20.55	30.65	45.82
	Energy	56.05	8.10	11.60	9.10	3.18	16.98	25.47
	ReAct	59.47	7.57	12.52	10.13	2.93	16.86	27.61
	MHE	17.59	9.20	7.68	4.74	0.33	8.99	15.86
	SHE	17.45	9.22	7.69	4.77	0.33	8.99	15.84
RecLag	18.12	6.40	4.60	2.67	0.28	6.82	12.09	
	± 2.02	± 0.25	± 0.12	± 0.47	± 0.02	± 0.13	± 0.25	± 0.01
ResNet34	MSP	59.86	28.26	32.06	31.69	33.61	43.28	45.56
	Energy	30.51	6.84	9.43	8.47	9.32	23.74	25.16
	ReAct	45.86	14.37	14.09	13.28	15.83	29.73	31.60
	MHE	6.20	6.17	4.40	2.94	2.34	14.32	15.86
	SHE	6.14	6.20	4.45	3.01	2.36	14.32	15.93
RecLag	5.19	5.60	2.85	2.11	2.31	12.04	11.71	
	± 0.24	± 0.07	± 0.05	± 0.05	± 0.03	± 0.07	± 0.23	± 0.11
WRN40-2	MSP	41.52	44.43	38.47	39.70	33.84	35.80	51.52
	Energy	15.35	17.77	14.98	17.45	10.58	19.71	36.75
	ReAct	18.83	19.93	18.25	20.68	11.98	21.67	42.02
	MHE	5.40	14.60	12.03	11.48	2.90	10.99	27.28
	SHE	5.25	14.39	13.18	12.39	2.83	10.98	28.35
RecLag	5.75	7.37	8.44	8.01	2.63	9.75	22.62	
	± 0.12	± 0.18	± 0.17	± 0.15	± 0.05	± 0.10	± 0.34	± 0.01

Objective

To evaluate OOD detection performance for RecLag.

Baselines

The FPR95 of the five OOD detection methods: {MSP, Energy, React, SHE, HE}

Settings

Features of a CIFAR-10 pretrained network are used for OOD detection. In-distribution data: CIFAR-10. OOD data: {SVHN, LSUN_C, LSUN_R, iSUN, Places, Tiny Image Net, SUN, iNaturalist}.

For the proposed RecLag the trimmed mean and standard deviations (following \pm symbols) of the largest and the smallest ones are reported.

Proposed Method

The same energy except the hidden-layer Lagrangian is replaced by following

$$\left(\frac{1}{\beta} \log \left(\frac{1}{\gamma} \sum_{\mu=1}^{N_H} e^{\beta \xi_\mu v} \right), 0 \right)$$

$$v^{k+1} = \Xi^T \text{softmax}(\beta \Xi v^k)$$

$$\log \left(\frac{1}{\gamma} \sum_{\mu=1}^{N_H} e^{\beta \xi_\mu v} \right)$$

$$(x \geq 0)$$

$$(x < 0)$$

$\gamma \in \mathbb{R}_{\geq 0}$: Real hyperparameter

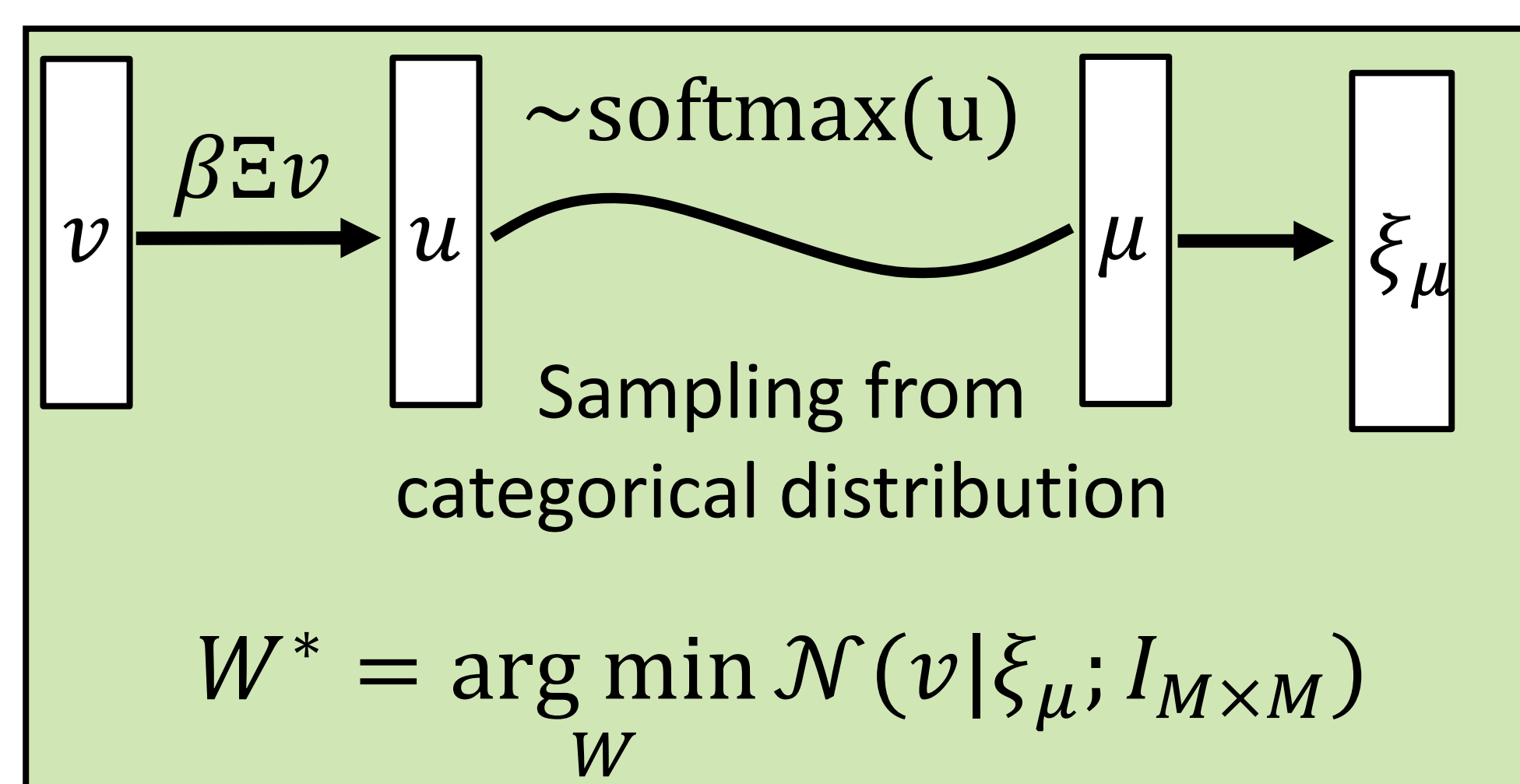
Probability-based learning

W is determined through optimization to satisfy the following probability equation.

$$p_H(\mu, v) = \frac{1}{Z} e^{\beta \xi_\mu v} \propto f(h)_\mu$$

$$\log(p(v)) - \log \left(\sum_{\mu=1}^{N_H} e^{\beta \xi_\mu v} \right) = \text{const}$$

$Z \in \mathbb{R}$: normalized coefficient

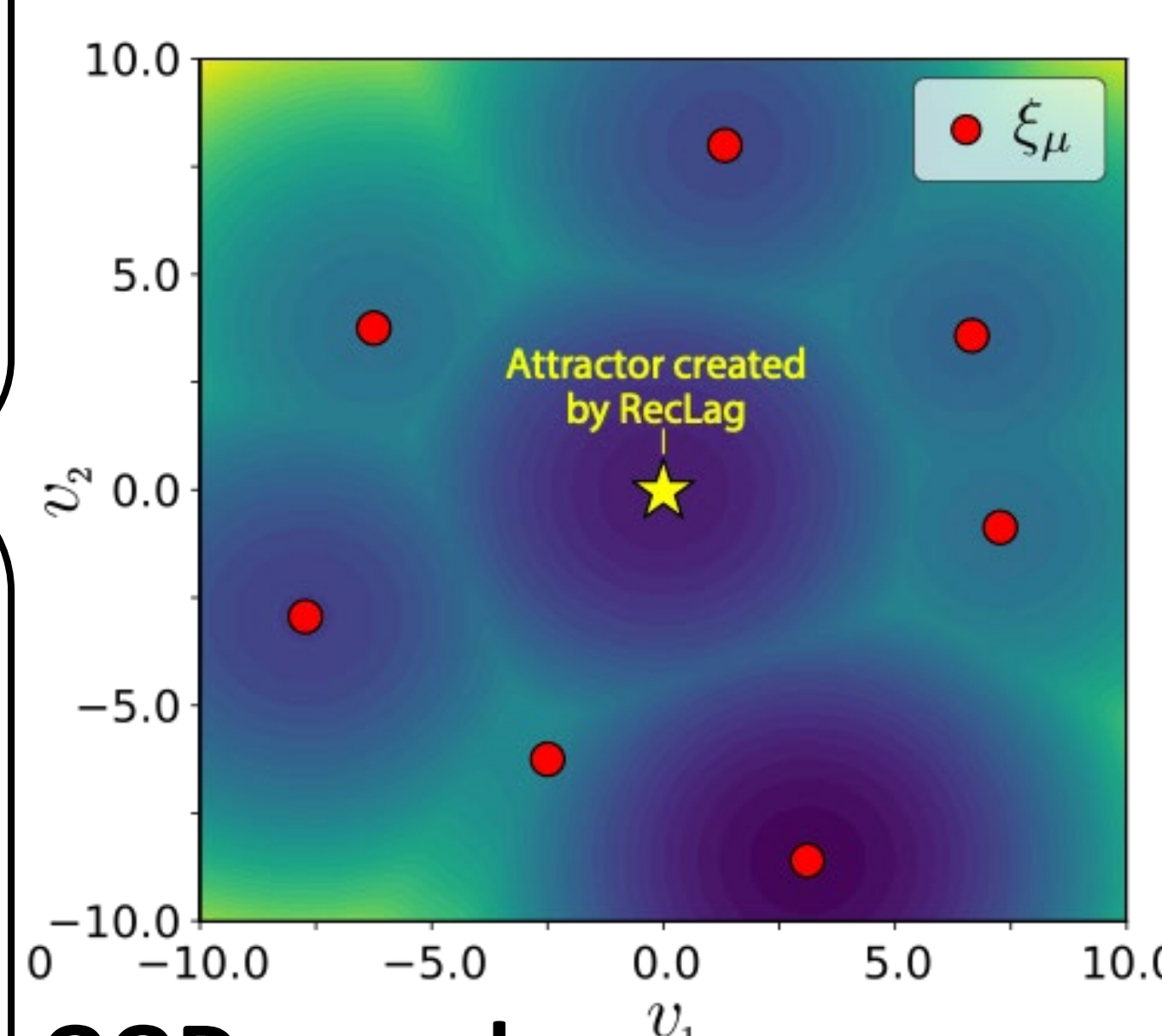


\mathcal{N} : multivariate normal distribution $I_{M \times M}$: M -dimensional identity matrix

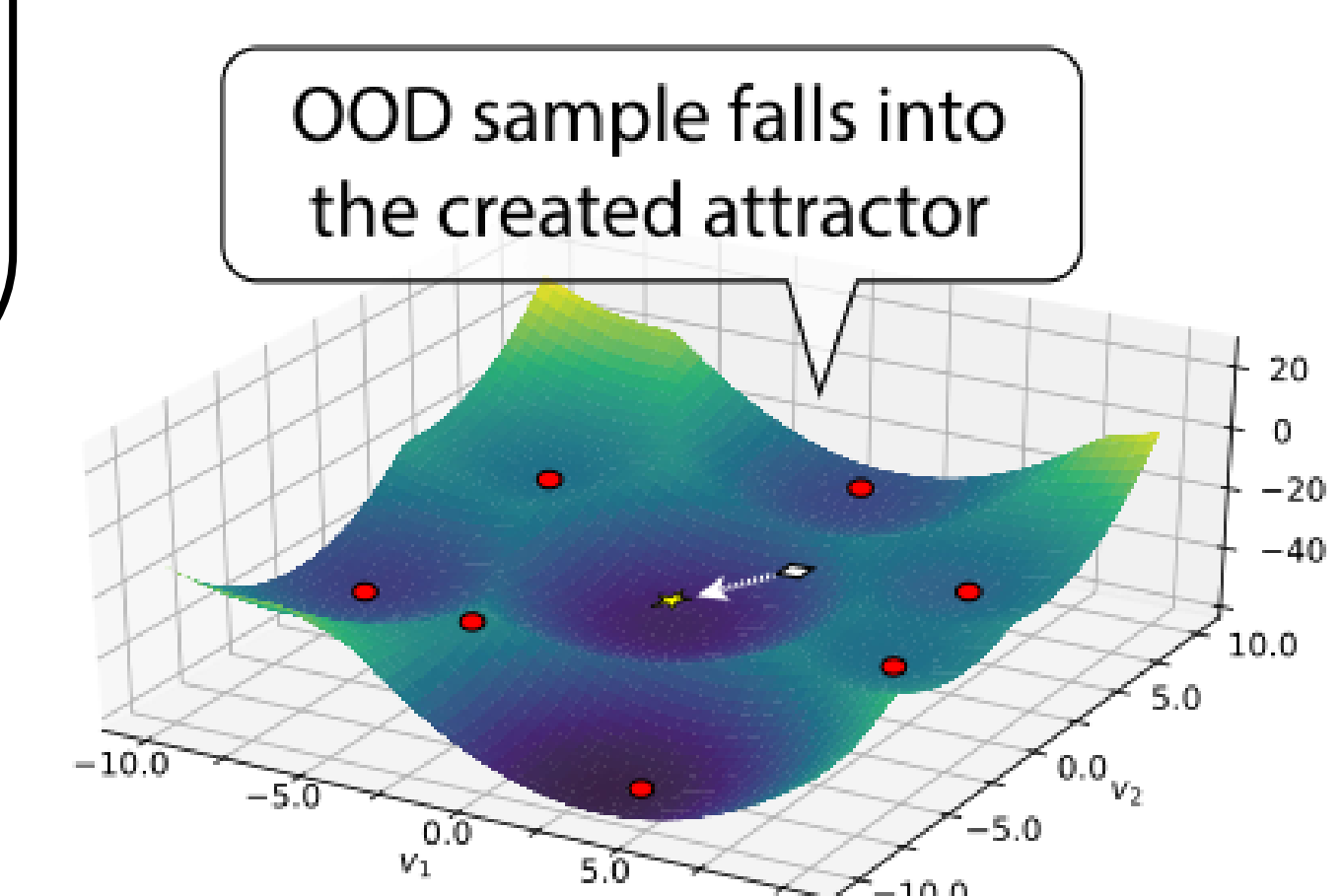
Each neuron in the hidden layer corresponds to the probability of the μ -th memory

The Lagrangian of the hidden layer corresponds to the probability density relative to the memory's parent distribution.

Energy Map



OOD sample



an out-of-distribution (OOD) query, MHN does retrieve a memory data, totally irrelevant.

Contributions

To address this, we propose the rectified Lagrangian (RegLag), a new Lagrangian for memory neurons that explicitly incorporates an attractor for OOD samples in the dynamical system of MHNs.

We demonstrate outperformance RecLag-based MHNs over energy-based OOD detection methods, including state-of-the-art Hopfield Energy, on nine image datasets.

Modern Hopfield Network (MHN)

Introduction [D. Krotov+ 2021, Widrich et al. 2021]

$$L = \sum_{i=1}^{N_V} v_i g_i(v) - L_V(v) + \sum_{\mu=1}^{N_H} h_\mu f_\mu(h) - L_H(h) - f(h)^T \Xi g(v)$$

Derivation functions

$$\frac{\partial L_H(h)}{\partial h}$$

$$g(v) = \frac{\partial L_V(v)}{\partial v}$$

Hidden-layer Lagrangian

$$\frac{1}{\beta} \log \left(\sum_{\mu=1}^{N_H} e^{\beta h_\mu} \right)$$

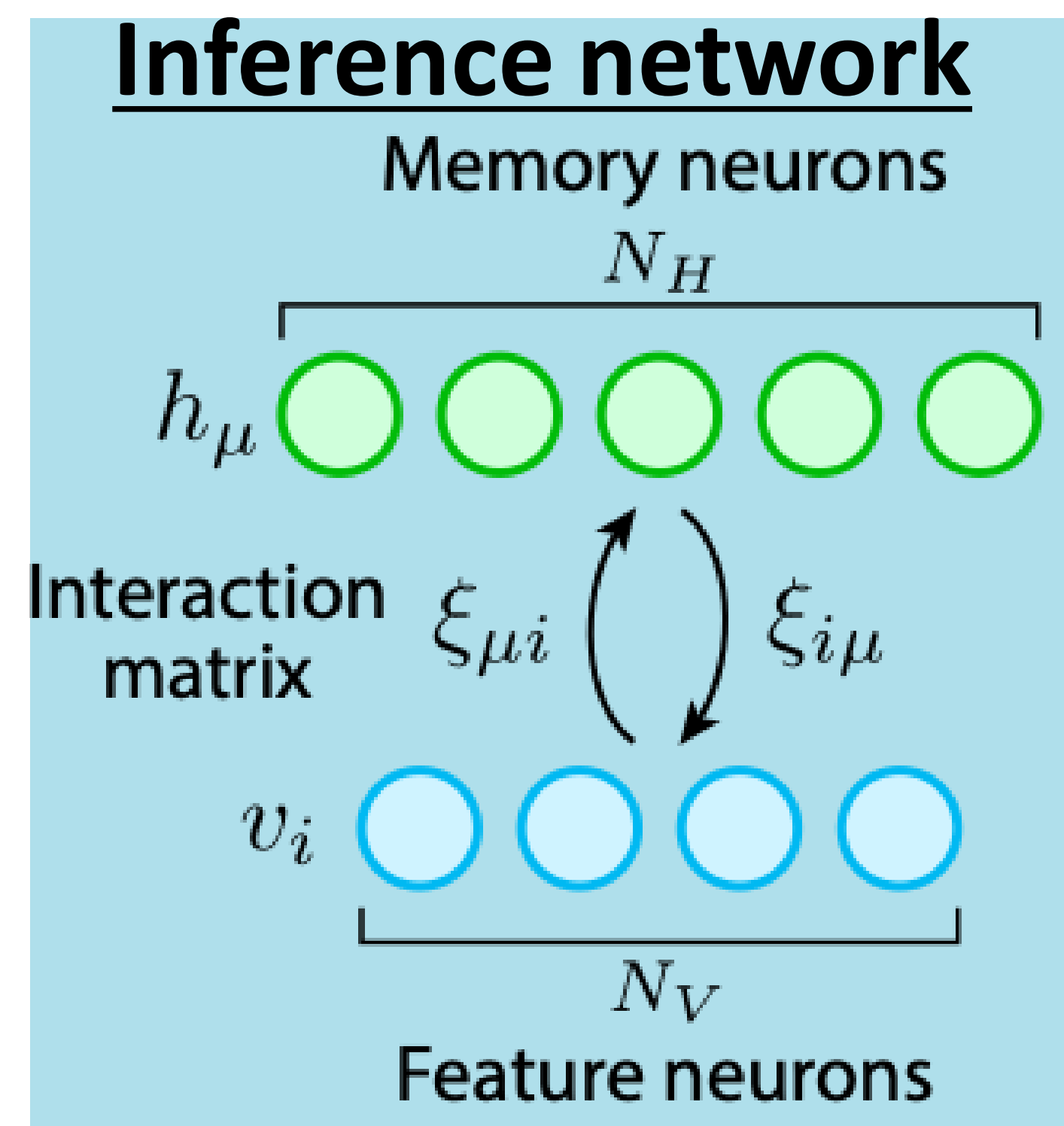
Visible Lagrangian

$$L_v = \frac{1}{2} \sum_{i=1}^{N_V} v_i^2$$

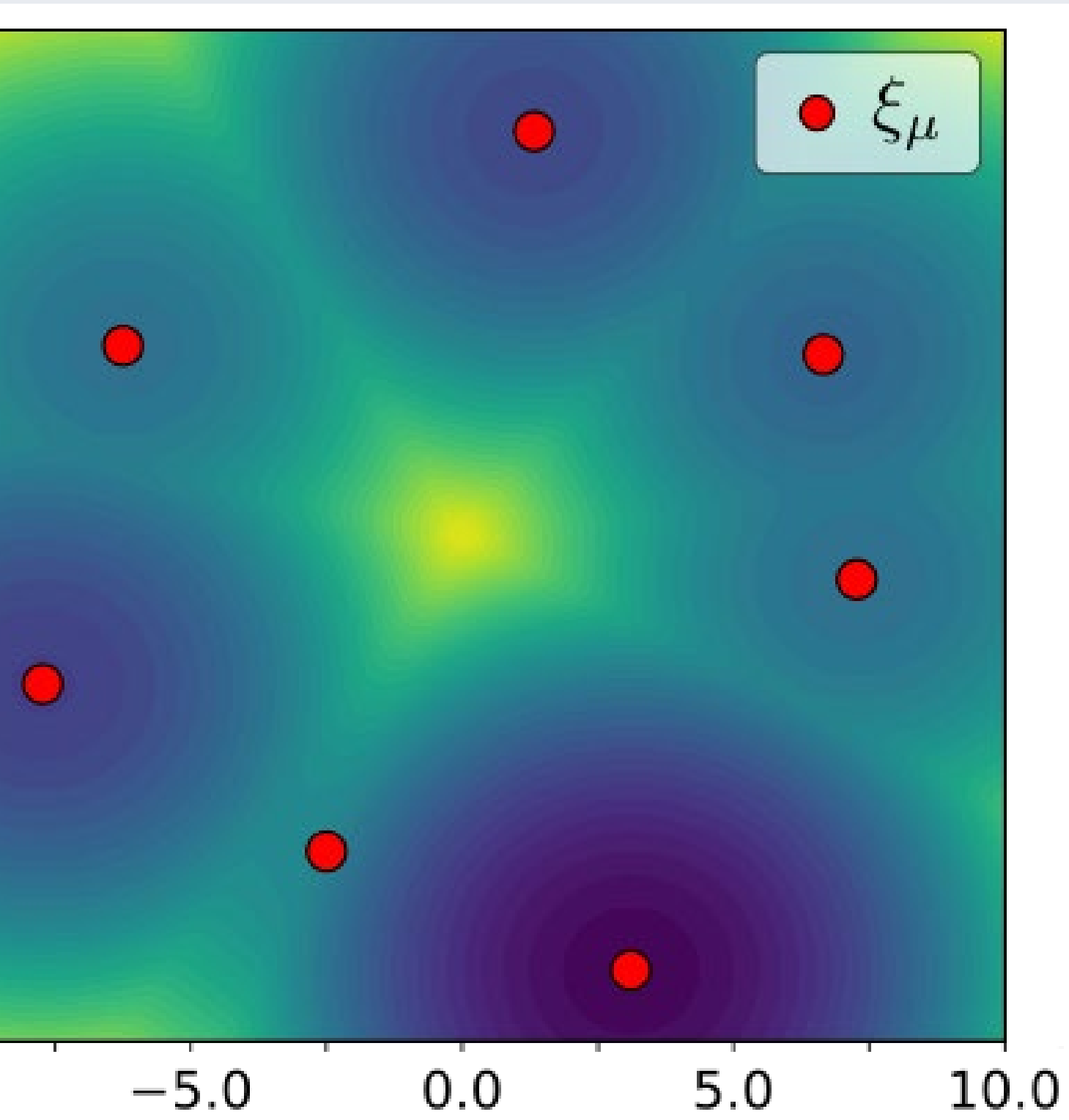
State Update rule

$$v^{k+1} = \Xi^T \text{softmax}(\beta \Xi v^k)$$

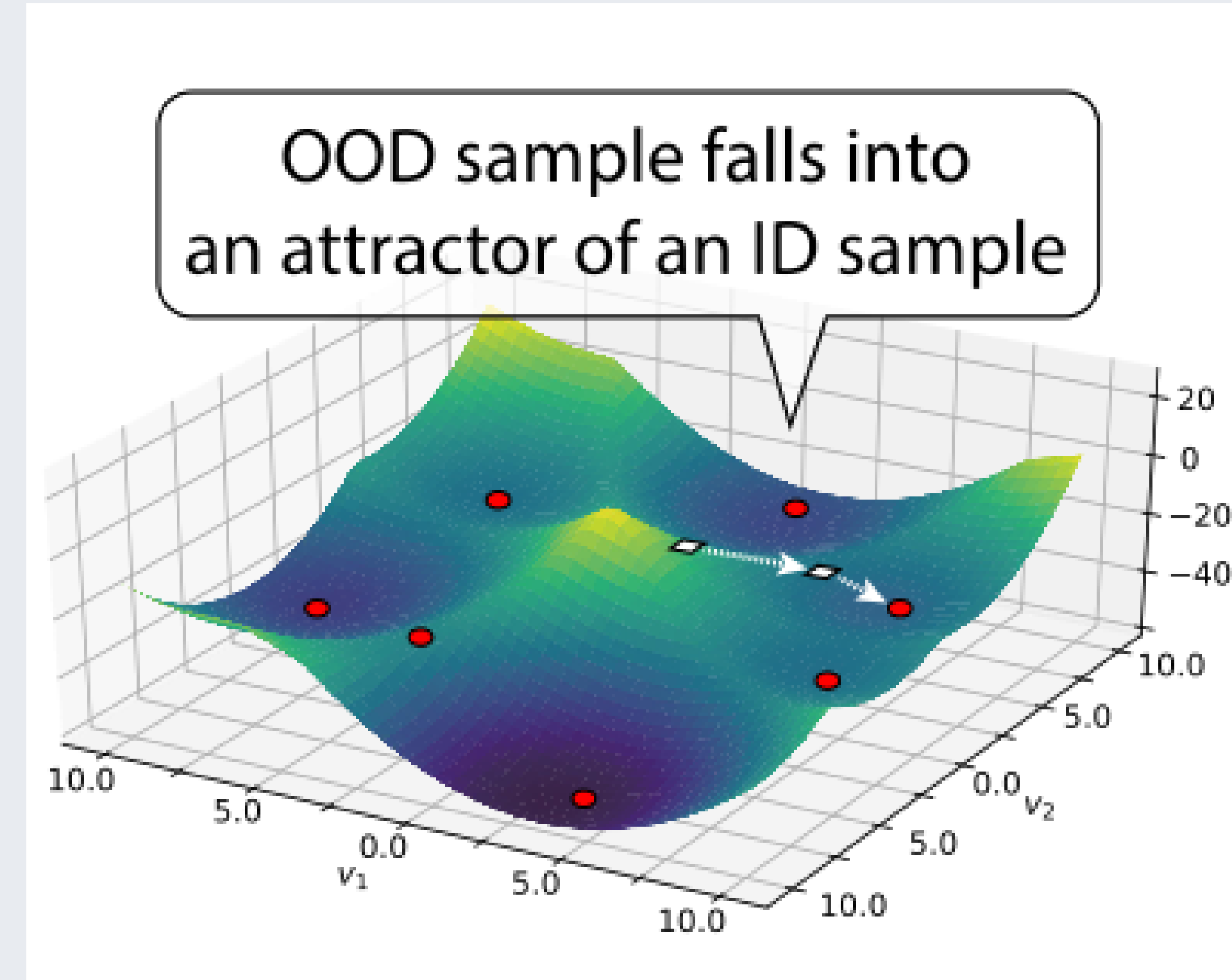
$v \in \mathbb{R}^{N_V}$: visible layer neuron $h \in \mathbb{R}^{N_H}$: hidden layer neuron, β : real hyperparameter, $\Xi \in \mathbb{R}^{N_H \times N_V}$: matrix parameter, ξ_μ : μ -th column vector component of a matrix Ξ



Energy Map



OOD sample



$$L_H(h) = \max \left(\frac{1}{\beta} \log \left(\frac{1}{\gamma} \sum_{\mu=1}^{N_H} e^{\beta h_\mu} \right), 0 \right)$$

Update rule

$$v^{k+1} = \chi \left(G(v^k) \right) \Xi^T \text{softmax}(\beta \Xi v^k)$$

$$G(v) = \frac{1}{\beta} \log \left(\frac{1}{\gamma} \sum_{\mu=1}^{N_H} e^{\beta \xi_\mu v} \right) \quad \chi(x) = \begin{cases} 1 & (x \geq 0) \\ 0 & (x < 0) \end{cases}$$

$\gamma \in \mathbb{R}_{\geq 0}$: Real hyperparameter

Probability-based learning

W is determined through optimization to satisfy the following probability equation

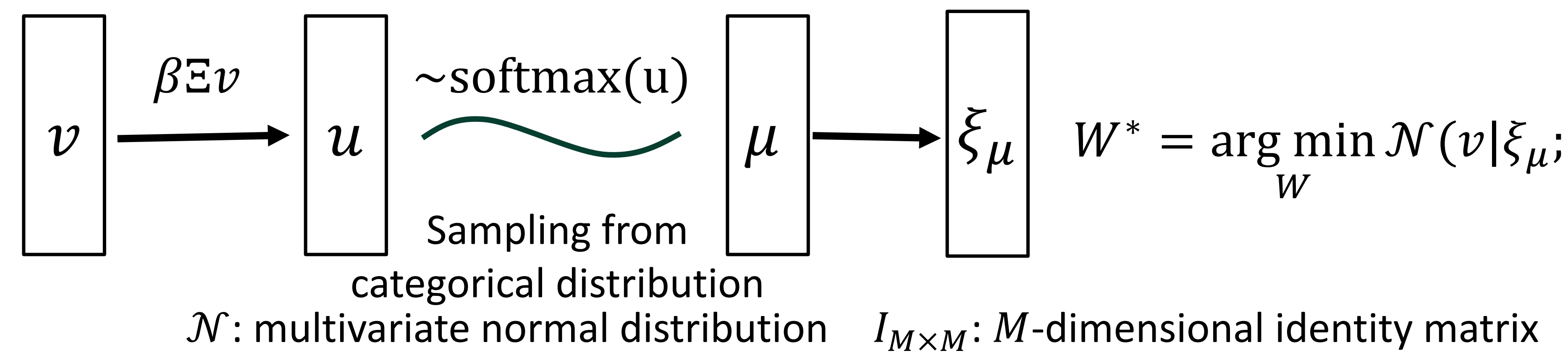
$$p_H(\mu, v) = \frac{1}{Z} e^{\beta \xi_\mu v} \propto f(h)_\mu$$

Each neuron in the hidden layer corresponds to the probability μ -th memory

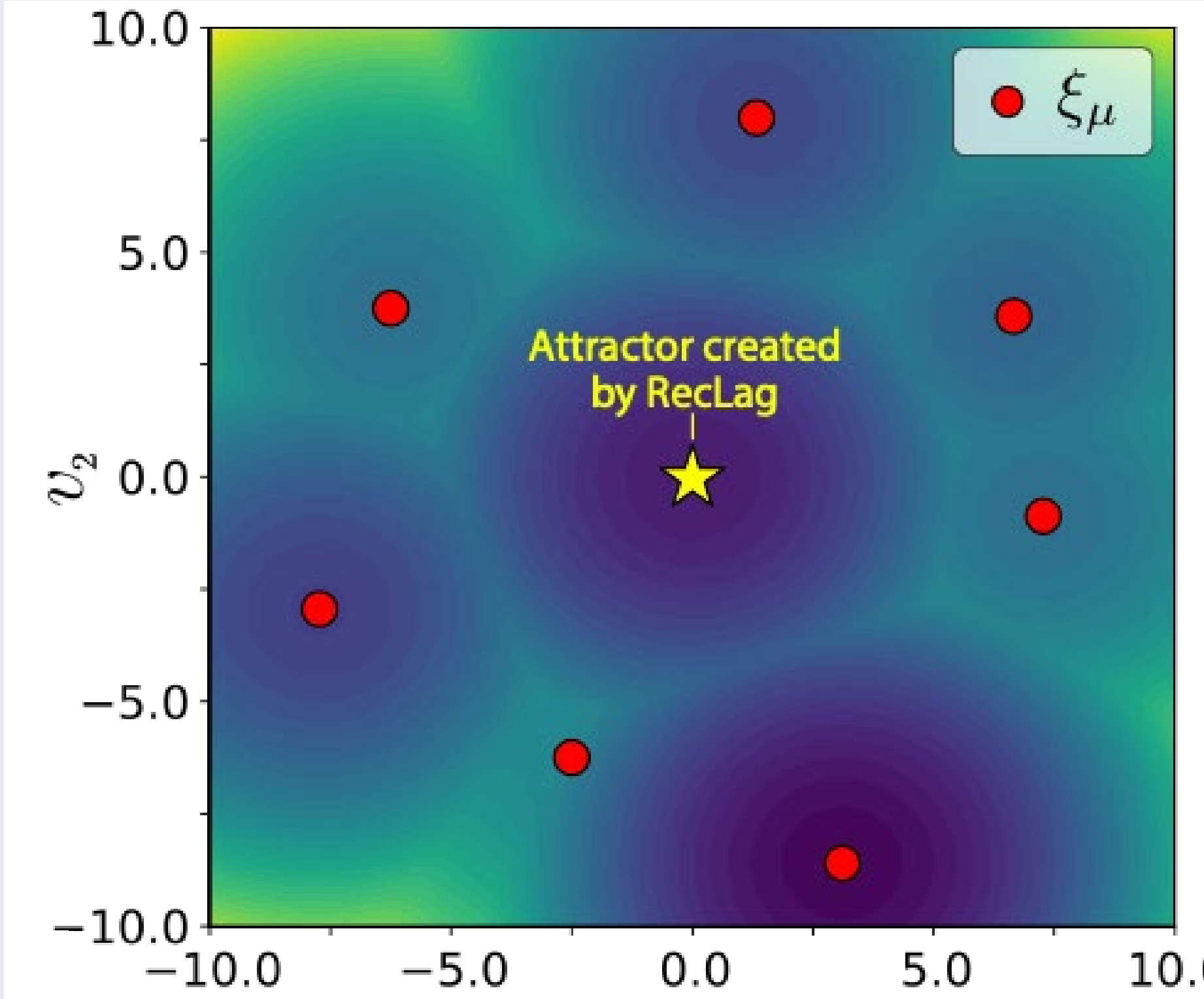
$$\log(p(v)) - \log \left(\sum_{\mu=1}^{N_H} e^{\beta \xi_\mu v} \right) = \text{const}$$

The Lagrangian of the hidden layer corresponds to the probability relative to the memory's parent distribution.

$Z \in \mathbb{R}$: normalized coefficient



Energy Map



OOD sample

