# Energy-based Hopfield Boosting for Out-of-Distribution Detection

Claus Hofmann [1 2]   Simon Schmid [3]   Bernhard Lehner [2 4]   Daniel Klotz [5]   Sepp Hochreiter [1 2]

## Abstract

Out-of-distribution (OOD) detection is critical when deploying machine learning models in the real world. Outlier exposure methods, which incorporate auxiliary outlier data in the training process, can drastically improve OOD detection performance compared to approaches without advanced training strategies. We introduce Hopfield Boosting, a boosting approach, which leverages modern Hopfield energy (MHE) to sharpen the decision boundary between the in-distribution and OOD data. Hopfield Boosting encourages the model to concentrate on hard-to-distinguish auxiliary outlier examples that lie close to the decision boundary between in-distribution and auxiliary outlier data. Our method achieves a new state-of-the-art in OOD detection with outlier exposure, improving the FPR95 metric from 2.28 to 0.92 on CIFAR-10 and from 11.76 to 7.94 on CIFAR-100.

## 1. Introduction

Out-of-distribution (OOD) detection is crucial when using machine learning systems in the real world. Deployed models will — sooner or later — encounter inputs that deviate from the training distribution. For example, a system trained to recognize music genres might also hear a sound clip of construction site noise. In the best case, a naive deployment can then result in overly confident predictions. In the worst case, we will get erratic model behavior and completely wrong predictions (Hendrycks & Gimpel, 2017). The purpose of OOD detection is to classify these examples as OOD, such that the system can then, for instance, notify users that no prediction is possible or trigger a manual investigation.

In this paper we propose Hopfield Boosting, a novel OOD detection method that leverages the energy component of modern Hopfield networks (MHNs; Ramsauer et al., 2021) and advances the state-of-the-art of OOD detection.

Our method uses an auxiliary outlier data set (AUX) to *boost* the model's OOD detection capacity. This allows the training process to learn a boundary around the in-distribution (ID) data, improving the performance at the OOD detection task. In summary, our contributions are as follows:

1. We propose Hopfield Boosting, an OOD detection approach that samples weak learners by using modern Hopfield energy (MHE; Ramsauer et al., 2021).

2. Hopfield Boosting achieves a new state-of-the-art in OOD detection. It improves the average false positive rate at 95% true positives (FPR95) from 2.28 to 0.92 on CIFAR-10 and from 11.38 to 7.94 on CIFAR-100.

3. We provide a theoretical background that motivates the suitability of Hopfield Boosting for OOD detection.

## 2. Related Work

Some authors (e.g., Bishop, 1994; Roth et al., 2022; Yang et al., 2022) distinguish between anomalies, outliers, and novelties. These distinctions reflect different goals within applications (Ruff et al., 2021). For example, when an anomaly is found, it will usually be removed from the training pipeline. However, when a novelty is found it should be studied. We focus on the detection of samples that are not part of the training distribution and consider sample categorization as a potential downstream task.

**Post-hoc OOD detection.** A common and straightforward OOD detection approach is to use a post-hoc strategy, where one employs statistics obtained from a classifier. The perhaps most well known and simplest approach in this class is the Maximum Softmax Probability (MSP; Hendrycks & Gimpel, 2017), where one utilizes $p(y \mid x)$ of the most likely class $y$ given a feature vector $x \in \mathbb{R}^D$ to estimate whether a sample is OOD. Despite good empirical performances, this view is intrinsically limited, since OOD detection should to focus on $p(x)$ (Morteza & Li, 2022). A

---

[1]Johannes Kepler University Linz, Institute for Machine Learning, JKU LIT SAL eSPML Lab, Austria [2]Silicon Austria Labs, JKU LIT SAL eSPML Lab, Austria [3]Software Competence Center Hagenberg GmbH, Austria [4]Johannes Kepler University Linz, JKU LIT SAL eSPML Lab, Austria [5]Department of Computational Hydrosystems, Helmholtz Centre for Environmental Research–UFZ, Leipzig, Germany. Correspondence to: Claus Hofmann <hofmann@ml.jku.at>.

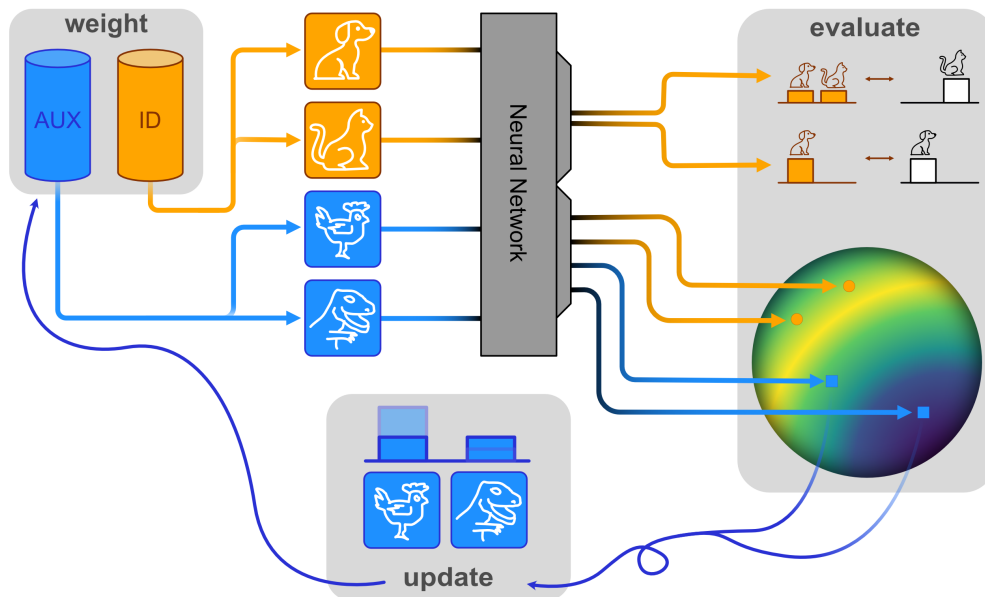Source code available at: https://github.com/ml-jku/hopfield-boosting

Figure 1: The Hopfield Boosting concept. The first step (weight) creates weak learners by firstly choosing in-distribution samples (ID, orange), and by secondly choosing auxiliary outlier samples (AUX, blue) according to their assigned probabilities; the second step (evaluate) computes the losses for the resulting predictions (Section 3); and the third step (update) assigns new probabilities to the AUX samples according to their position on the hypersphere (see Figure 2).

wide range of post-hoc OOD detection approaches have been proposed to address the shortcomings of MSP (e.g., Lee et al., 2018; Liang et al., 2018; Hendrycks et al., 2019a; Liu et al., 2020; Huang et al., 2021; Sun et al., 2021; 2022; Wang et al., 2022; Sun et al., 2022; Sun & Li, 2022; Zhang et al., 2023a; Djurisic et al., 2023; Liu et al., 2023; Xu et al., 2024). Most related to Hopfield Boosting is the work of Zhang et al. (2023a), who, to our knowledge, are the first to apply MHE for OOD detection. Specifically, they propose a post-hoc OOD detection approach using MHE. They use the ID data set to produce stored patterns and then use a modified version of MHE as the OOD score. While post-hoc approaches can be deployed out of the box on any model, a crucial limitation is that their performance heavily depends on the employed model itself.

**Training methods.** In contrast to post-hoc strategies, training-based methods modify the training process to improve the model's OOD detection capability (e.g., Hendrycks et al., 2019c; Tack et al., 2020; Sehwag et al., 2021; Du et al., 2022; Hendrycks et al., 2022; Wei et al., 2022a; Ming et al., 2023; Tao et al., 2023; Lu et al., 2024). For example, Self-Supervised Outlier Detection (SSD; Sehwag et al., 2021) leverages contrastive self-supervised learning to train a model for OOD detection.

**Auxiliary outlier data and outlier exposure.** A third group of OOD detection approaches are outlier exposure

(OE) methods. Like Hopfield Boosting, they incorporate AUX data in the training process to improve the detection of OOD samples (e.g., Hendrycks et al., 2019b; Liu et al., 2020; Ming et al., 2022). As far as we know, all OE approaches optimize an objective ($\mathcal{L}_{\text{OOD}}$), which aims at improving the model's discriminative power between ID and OOD data using the AUX data set as a stand-in for the OOD case. Hendrycks et al. (2019b) were the first to use the term OE to describe a more restrictive OE concept. Since their approach uses the MSP for incorporating the AUX data we refer to it as MSP-OE. Further, we refer to the OE approach introduced in Liu et al. (2020) as EBO-OE (to differentiate it from EBO, their post-hoc approach). In general, OE methods conceptualize the AUX data set as a large and diverse data set (e.g., ImageNet for vision tasks). As a consequence, usually, only a small subset of the samples bear semantic similarity to the ID data set — most data points are easily distinguishable from the ID data set. Recent approaches therefore actively try to find informative samples for the training. The aim is to refine the decision boundary, ensuring the ID data is more tightly encapsulated (e.g., Chen et al., 2021; Ming et al., 2022; Jiang et al., 2024). For example, Posterior Sampling-based Outlier Mining (POEM; Ming et al., 2022) selects samples close to the decision boundary using Thompson sampling: They first sample a linear decision boundary between ID and AUX data and then select those data instances which are closest to the sampled decision boundary. Hopfield Boosting also makes use of samples close to the

boundary by giving them higher weights for the boosting step. Further OE methods include mixing-based OE methods (Zhang et al., 2023b; Zhu et al., 2023a), that employ mixup (Zhang et al., 2018) or CutMix (Yun et al., 2019) to augment the OE task. Distribution-Augmented OOD Learning (DAL; Wang et al., 2023a) augments the AUX data by defining a Wasserstein-1 ball around the AUX data and performing OE using this Wasserstein ball. Diversified Outlier Exposure (DivOE; Zhu et al., 2023b) synthesizes artificial auxiliary outliers, resulting in a more diverse coverage of the feature space. Distributional-agnostic Outlier Exposure (DOE; Wang et al., 2023b) implicitly synthesizes auxiliary outlier data using a transformation of the model weights, arguing that perturbing model parameters has the same effect as transforming the data. Diverse Outlier Sampling (DOS; Jiang et al., 2024) applies K-Means clustering to the features of the AUX data set and then samples outlier points close to the decision boundary uniformly from the obtained clusters. All mentioned OE-based methods were able to outperform their respective contemporary post-hoc or training approaches by a sizeable margin. This indicates that incorporating an appropriate AUX data set and applying OE during the training process improves the OOD detection capability in general.

**Continuous modern Hopfield networks.** MHNs are energy-based associative memory networks. They advance conventional Hopfield networks (Hopfield, 1984) by introducing continuous queries and states with the MHE as a new energy function. MHE leads to exponential storage capacity, while retrieval is possible with a one-step update (Ramsauer et al., 2021). The update rule of MHNs coincides with the softmax-attention as it is used in the Transformer (Vaswani et al., 2017). Examples for successful applications of MHNs are Widrich et al. (2020); Fürst et al. (2022); Sanchez-Fernandez et al. (2022); Paischer et al. (2022); Schäfl et al. (2022) and Auer et al. (2023). Section 3.2 gives an introduction to MHE for OOD detection. For further details on MHNs, we refer to Appendix A.

**Boosting for classification.** Boosting, in particular, AdaBoost (Freund & Schapire, 1995), is an ensemble learning technique for classification. It is designed to focus ensemble members toward data instances that are hard to classify by assigning them higher weights. These challenging instances often lie near the maximum margin hyperplane (Rätsch et al., 2001), akin to support vectors in support vector machines (SVMs; Cortes & Vapnik, 1995). Popular boosting methods include Gradient boosting (Breiman, 1997), LogitBoost (Friedman et al., 2000), and LPBoost (Demiriz et al., 2002).

**Radial basis function networks.** Radial basis function networks (RBF networks; Moody & Darken, 1989) are func-

tion approximators of the form

$$\varphi(\boldsymbol{\xi}) = \sum_{i=1}^{N} \omega_i \exp\left(-\frac{||\boldsymbol{\xi} - \boldsymbol{\mu}_i||_2^2}{2\sigma_i^2}\right), \quad (1)$$

where $\omega_i$ are linear weights, $\boldsymbol{\mu}_i$ are the component means and $\sigma_i^2$ are the component variances. RBF networks can be described as a weighted linear superposition of $N$ radial basis functions and have previously been used as hypotheses for boosting (Rätsch et al., 2001). If the linear weights are strictly positive, RBF networks can be viewed as an unnormalized weighted mixture of Gaussian distributions $p_i(\boldsymbol{\xi}) = \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}_i, \sigma_i^2 \boldsymbol{I})$ with $i = \{1, \dots, N\}$. Appendix E.1 explores the connection between RBF networks and MHNs via Gaussian mixtures in more depth. We refer to Bishop (1995) and Müller et al. (1997) for more general information on RBF networks.

## 3. Method

This section presents Hopfield Boosting: First, we formalize the OOD detection task. Second, we give an overview of MHE and why it is suitable for OOD detection. Finally, we introduce the AUX-based boosting framework. Figure 1 shows a summary of the Hopfield Boosting concept.

### 3.1. Classification and OOD Detection

Consider a multi-class classification problem denoted as $(\boldsymbol{X}^{\mathcal{D}}, \boldsymbol{Y}^{\mathcal{D}}, \mathcal{Y})$, where $\boldsymbol{X}^{\mathcal{D}} \in \mathbb{R}^{D \times N}$ represents a set of $N$ $D$-dimensional feature vectors $(\boldsymbol{x}_1^{\mathcal{D}}, \boldsymbol{x}_2^{\mathcal{D}}, \dots, \boldsymbol{x}_N^{\mathcal{D}})$, which are i.i.d. samples $\boldsymbol{x}_i^{\mathcal{D}} \sim p_{\text{ID}}$. $\boldsymbol{Y}^{\mathcal{D}} \in \mathcal{Y}^N$ corresponds to the labels associated with these feature vectors, and $\mathcal{Y}$ is a set containing possible classes ($||\mathcal{Y}|| = K$ signifies the number of distinct classes present).

We consider observations $\boldsymbol{\xi}^{\mathcal{D}} \in \mathbb{R}^D$ that deviate considerably from the data generation $p_{\text{ID}}(\boldsymbol{\xi}^{\mathcal{D}})$ that defines the "normality" of our data as OOD. Following Ruff et al. (2021), an observation is OOD if it pertains to the set

$$\mathbb{O} = \{\boldsymbol{\xi}^{\mathcal{D}} \in \mathbb{R}^D \mid p_{\text{ID}}(\boldsymbol{\xi}^{\mathcal{D}}) < \epsilon\} \text{ where } \epsilon \geq 0. \quad (2)$$

Since the probability density of the data generation $p_{\text{ID}}$ is in general not known, one needs to estimate $p_{\text{ID}}(\boldsymbol{\xi}^{\mathcal{D}})$. In practice, it is common to define an outlier score $s(\boldsymbol{\xi})$ that uses an encoder $\phi$, where $\boldsymbol{\xi} = \phi(\boldsymbol{\xi}^{\mathcal{D}})$. The outlier score should — in the best case — preserve the density ranking.

In contrast to a density estimation, the score $s(\boldsymbol{\xi})$ does not have to fulfill all requirements of a probability density (like proper normalization or non-negativity). Given $s(\boldsymbol{\xi})$ and $\phi$, OOD detection can be formulated as a binary classification task with the classes ID and OOD:

$$\hat{B}(\boldsymbol{\xi}^{\mathcal{D}}, \gamma) = \begin{cases} \text{ID} & \text{if } s(\phi(\boldsymbol{\xi}^{\mathcal{D}})) \geq \gamma \\ \text{OOD} & \text{if } s(\phi(\boldsymbol{\xi}^{\mathcal{D}})) < \gamma \end{cases}. \quad (3)$$

It is common to choose the threshold $\gamma$ so that a portion of 95% of ID samples from a previously unseen validation set are correctly classified as ID. However, metrics like the area under the receiver operating characteristic (AUROC) can be directly computed on $s(\boldsymbol{\xi})$ without specifying $\gamma$ since the AUROC computation sweeps over the threshold.

## 3.2. Modern Hopfield Energy

The log-sum-exponential (lse) is defined as

$$\mathrm{lse}(\beta, \boldsymbol{z}) = \beta^{-1} \log \left( \sum_{i=1}^{N} \exp(\beta z_i) \right), \qquad (4)$$

where $\beta$ is the inverse temperature and $\boldsymbol{z} \in \mathbb{R}^N$ is a vector. The lse can be seen as a soft approximation to the maximum function: As $\beta \to \infty$, the lse approaches $\max_i z_i$.

Given a set of $N$ $d$-dimensional stored patterns $(\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N)$ arranged in a data matrix $\boldsymbol{X}$, and a $d$-dimensional query $\boldsymbol{\xi}$, the MHE is defined as

$$\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X}) = -\mathrm{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) + \frac{1}{2} \boldsymbol{\xi}^T \boldsymbol{\xi} + C, \quad (5)$$

where $C = \beta^{-1} \log N + \frac{1}{2} M^2$ and $M$ is the largest norm of a pattern: $M = \max_i \|x_i\|$. For a fixed data matrix $\boldsymbol{X}$, $C$ remains constant for any given query $\boldsymbol{\xi}$. $\boldsymbol{X}$ is also called the memory of the MHN. Intuitively, Equation (5) can be explained as follows: The dot-product within the lse computes a similarity for a given $\boldsymbol{\xi} \in \mathbb{R}^d$ to all patterns in the memory $\boldsymbol{X} \in \mathbb{R}^{d \times N}$. The lse function aggregates the similarities to form a single value. The parameter $\beta$ controls how the similarities are aggregated. If $\beta \to \infty$, the similarity of $\boldsymbol{\xi}$ to the closest pattern in $\boldsymbol{X}$ is returned. Tuning $\beta$ allows one to adjust the "softness" of the lse. We can regard $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ as an expression of dissimilarity between $\boldsymbol{X}$ and $\boldsymbol{\xi}$: If $\boldsymbol{\xi}$ is similar to $\boldsymbol{X}$, $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ will be low. In contrast, if $\boldsymbol{\xi}$ is dissimilar from $\boldsymbol{X}$, $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ will be high.

To use MHE for OOD detection, one needs a memory $\boldsymbol{X}$ of patterns from the ID data set. Hopfield Boosting acquires the memory patterns in $\boldsymbol{X}$ by feeding raw data instances $(\boldsymbol{x}_1^{\mathcal{D}}, \boldsymbol{x}_2^{\mathcal{D}}, \ldots, \boldsymbol{x}_N^{\mathcal{D}})$ arranged in the data matrix $\boldsymbol{X}^{\mathcal{D}}$ of the ID data set to an encoder $\phi : \mathbb{R}^D \to \mathbb{R}^d$ (e.g., ResNet): $\boldsymbol{x}_i = \phi(\boldsymbol{x}_i^{\mathcal{D}})$. We denote this component-wise application of an encoder to obtain a memory matrix as $\boldsymbol{X} = \phi(\boldsymbol{X}^{\mathcal{D}})$. Similarly, a raw query $\boldsymbol{\xi}^{\mathcal{D}}$ is fed through the encoder to obtain the query pattern: $\boldsymbol{\xi} = \phi(\boldsymbol{\xi}^{\mathcal{D}})$. One can now use $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ to estimate whether $\boldsymbol{\xi}$ is ID or OOD: A low energy indicates $\boldsymbol{\xi}$ is ID, and a high energy signifies that $\boldsymbol{\xi}$ is OOD.

## 3.3. Boosting Framework

**Sampling of informative outlier data.** Hopfield Boosting uses an AUX data set to learn a decision boundary between the ID and OOD region during the training process. Similar to Chen et al. (2021) and Ming et al. (2022), Hopfield Boosting selects informative outliers close to the ID-OOD decision boundary and uses them for training the model. For this selection, Hopfield Boosting weights the AUX data similar to AdaBoost (Freund & Schapire, 1995) by sampling data instances close to the decision boundary more frequently. We consider samples close to the decision boundary as weak learners — their nearest neighbors consist of samples from their own class as well as from the foreign class. An individual weak learner would yield a classifier that is only slightly better than random guessing (Figure 7). Vice versa, a strong learner can be created by forming an ensemble of a set of weak learners (Figure 2).

We denote the matrix containing the raw AUX data instances $(\boldsymbol{o}_1^{\mathcal{D}}, \boldsymbol{o}_2^{\mathcal{D}}, \ldots, \boldsymbol{o}_N^{\mathcal{D}})$ as $\boldsymbol{O}^{\mathcal{D}} \in \mathbb{R}^{D \times M}$, and the memory containing the encoded AUX patterns as $\boldsymbol{O} = \phi(\boldsymbol{O}^{\mathcal{D}})$. The boosting process proceeds as follows: There exists a weight $(w_1, w_2, \ldots, w_N)$ for each data point in $\boldsymbol{O}^{\mathcal{D}}$ and the individual weights are aggregated into the weight vector $\boldsymbol{w}_t$. Hopfield Boosting uses these weights to draw mini-batches $\boldsymbol{O}_s^{\mathcal{D}}$ from $\boldsymbol{O}^{\mathcal{D}}$ for training, where weak learners are sampled more frequently into the Hopfield memory.

We introduce an MHE-based energy function which Hopfield Boosting uses to determine how weak a specific learner $\boldsymbol{\xi}$ is (with higher energy indicating a weaker learner):

$$\begin{aligned} \mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) = &-2\,\mathrm{lse}(\beta, (\boldsymbol{X} \,\|\, \boldsymbol{O})^T \boldsymbol{\xi}) + \\ &\mathrm{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) + \mathrm{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}), \end{aligned} \qquad (6)$$

where $\boldsymbol{X} \in \mathbb{R}^{d \times N}$ contains ID patterns, $\boldsymbol{O} \in \mathbb{R}^{d \times M}$ contains AUX patterns, and $(\boldsymbol{X} \,\|\, \boldsymbol{O}) \in \mathbb{R}^{d \times (N+M)}$ denotes the concatenated data matrix containing the patterns from both $\boldsymbol{X}$ and $\boldsymbol{O}$. Before computing $\mathrm{E}_b$, we normalize the feature vectors in $\boldsymbol{X}, \boldsymbol{O}$, and $\boldsymbol{\xi}$ to unit length.

Figure 3 displays the energy landscape of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ using exemplary data on a 3-dimensional sphere. $\mathrm{E}_b$ is maximal at the decision boundary between ID and AUX data and decreases with increasing distance from the decision boundary in both directions.

As we show in our theoretical discussion in Appendix D, when modeling the class-conditional densities of the ID and AUX data set as mixtures of Gaussian distributions

$$p(\boldsymbol{\xi} \mid \mathrm{ID}) = \frac{1}{N} \sum_{i=1}^{N} \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{x}_i, \beta^{-1} \boldsymbol{I}), \qquad (7)$$

$$p(\boldsymbol{\xi} \mid \mathrm{AUX}) = \frac{1}{N} \sum_{i=1}^{N} \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{o}_i, \beta^{-1} \boldsymbol{I}), \qquad (8)$$

with equal class priors $p(\mathrm{ID}) = p(\mathrm{AUX}) = 1/2$ and normalized patterns $\|\boldsymbol{x}_i\| = 1$ and $\|\boldsymbol{o}_i\| = 1$, we obtain

$E_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) \overset{C}{=} \beta^{-1} \log(p(\text{ ID} \mid \boldsymbol{\xi}) \cdot p(\text{ AUX} \mid \boldsymbol{\xi}))$, where $\overset{C}{=}$ denotes equality up to an irrelevant additive constant. The exponential of $E_b$ is the variance of a Bernoulli random variable with the outcomes $\{\text{ID}, \text{AUX}\}$ conditioned on $\boldsymbol{\xi}$. Thus, according to $E_b$, the weak learners are situated at exactly those locations in $\mathbb{R}^d$ where the variance of the model's decision between the ID and AUX classes is high.

Given a set of query values $(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \ldots, \boldsymbol{\xi}_n)$ assembled in a query matrix $\boldsymbol{\Xi} \in \mathbb{R}^{d \times n}$, we denote a vector of energies $\boldsymbol{e} \in \mathbb{R}^n$ with $e_i = E_b(\boldsymbol{\xi}_i; \boldsymbol{X}, \boldsymbol{O})$ as

$$\boldsymbol{e} = E_b(\boldsymbol{\Xi}; \boldsymbol{X}, \boldsymbol{O}). \tag{9}$$

To calculate the weights $\boldsymbol{w}_{t+1}$, we use the memory of AUX patterns as a query matrix $\boldsymbol{\Xi} = \boldsymbol{O}$ and compute the respective energies $E_b$ of those patterns. The resulting energy vector $E_b(\boldsymbol{\Xi}; \boldsymbol{X}, \boldsymbol{O})$ is then normalized by a $\mathrm{softmax}$. This computation provides the updated weights:

$$\boldsymbol{w}_{t+1} = \mathrm{softmax}(\beta E_b(\boldsymbol{\Xi}; \boldsymbol{X}, \boldsymbol{O})). \tag{10}$$

Appendix G provides theoretical background on how informative samples close to the decision boundary are beneficial for training an OOD detector.

**Training the model with MHE.** In this section, we introduce how Hopfield Boosting uses the sampled weak learners to improve the detection of patterns outside the training distribution. We follow the established training method for OE (Hendrycks et al., 2019b; Liu et al., 2020; Ming et al., 2022): Train a classifier on the ID data using the standard cross-entropy loss and add an OOD loss that uses the AUX data set to sharpen the decision boundary between the ID and OOD regions. Formally, this yields a composite loss

$$\mathcal{L} = \mathcal{L}_{\mathrm{CE}} + \lambda \mathcal{L}_{\mathrm{OOD}}, \tag{11}$$

where $\lambda$ is a hyperparamter indicating the relative importance of $\mathcal{L}_{\mathrm{OOD}}$. Hopfield Boosting explicitly minimizes $E_b$ (which is also the energy function Hopfield Boosting uses to sample weak learners). Given the weight vector $\boldsymbol{w}_t$, and the data sets $\boldsymbol{X}^{\mathcal{D}}$ and $\boldsymbol{O}^{\mathcal{D}}$, we obtain a mini-batch $\boldsymbol{X}_s^{\mathcal{D}}$ containing N samples from $\boldsymbol{X}^{\mathcal{D}}$ by uniform sampling, and a mini-batch of N weak learners $\boldsymbol{O}_s^{\mathcal{D}}$ from $\boldsymbol{O}^{\mathcal{D}}$ by sampling according to $\boldsymbol{w}_t$ with replacement. We then feed the respective mini-batches into the neural network $\phi_{\mathrm{base}}$ to create a latent representation (in our experiments, we always use the representation of the penultimate layer of a ResNet). Our proposed approach then uses two heads:

1. A linear classification head that maps the latent representation to the class logits for $\mathcal{L}_{\mathrm{CE}}$.

2. A 2-layer MLP $\phi_{\mathrm{proj}}$ maps the representations from the penultimate layer to the output for $\mathcal{L}_{\mathrm{OOD}}$.

Hopfield Boosting computes $\mathcal{L}_{\mathrm{OOD}}$ on the representations it obtains from $\phi = \phi_{\mathrm{proj}} \circ \phi_{\mathrm{base}}$ as follows:

$$\mathcal{L}_{\mathrm{OOD}} = \frac{1}{2N} \sum_{\boldsymbol{\xi}} E_b(\boldsymbol{\xi}; \boldsymbol{X}_s, \boldsymbol{O}_s), \tag{12}$$

where the memories $\boldsymbol{X}_s$ and $\boldsymbol{O}_s$ contain the encodings of the sampled data instances: $\boldsymbol{X}_s = \phi(\boldsymbol{X}_s^{\mathcal{D}})$ and $\boldsymbol{O}_s = \phi(\boldsymbol{O}_s^{\mathcal{D}})$. The sum is taken over the observations $\boldsymbol{\xi}$, which are drawn from $(\boldsymbol{X}_s \parallel \boldsymbol{O}_s)$. Hopfield Boosting computes $\mathcal{L}_{\mathrm{OOD}}$ for each mini-batch by first calculating the pairwise similarity matrix between the patterns in the mini-batch, followed by determining the $E_b$ values of the individual observations $\boldsymbol{\xi}$, and, finally a mean reduction. To the best of our knowledge, Hopfield Boosting is the first method that uses Hopfield networks in this way to train a deep neural network. We note that there is a relation between Hopfield Boosting and SVMs with an RBF kernel (see Appendix E.3). However, the optimization procedure of SVMs is in general not differentiable. In contrast, our novel energy function is fully differentiable. This allows us to use it to train neural networks.

---

**Algorithm 1** Hopfield Boosting

**Require:** $T, N, \boldsymbol{X}, \boldsymbol{O}, \boldsymbol{Y}, \mathcal{L}_{\mathrm{CE}}, E_b, \beta$
  Set all weights $w_1$ to $1/|\boldsymbol{O}|$
  **for** $t = 1$ to $T$ **do**
    1. **Weight**. Get hypothesis $\boldsymbol{X}_s \parallel \boldsymbol{O}_s \to \{\text{ID}, \text{AUX}\}$:
      1.a. Mini-batch sampling $\boldsymbol{X}_s$ from $\boldsymbol{X}$, and
      1.b. Sub-sampling of weak learners $\boldsymbol{O}_s$ from $\boldsymbol{O}$
        according to the weighting $\boldsymbol{w}_t$.
    2. **Evaluate**. Compute composite loss from
      Equation (11) on $\boldsymbol{X}_s$ and $\boldsymbol{O}_s$.
    3. **Update**. Update model for the next iteration:
      3.a. At every step, update the full model
        (backbone, classification head, and MHE).
      3.b. At every $t * N$ step calculate new weights for
        $\boldsymbol{O}$ with $\boldsymbol{w}_{t+1} = \mathrm{softmax}(\beta E_b(\boldsymbol{O}; \boldsymbol{X}, \boldsymbol{O}))$.
  **end for**
return $\boldsymbol{w}_t$

---

**Summary.** Algorithm 1 provides an outline of Hopfield Boosting. Each iteration $t$ consists of three main steps: 1. weight, 2. evaluate, and 3. update. In the first step, Hopfield Boosting samples a mini-batch from the ID data and **weights** the AUX data by sampling a mini-batch according to $\boldsymbol{w}_t$. Then, Hopfield Boosting **evaluates** the composite loss on the sampled mini-batch. Lastly, Hopfield Boosting **updates** the model parameters and, every $N$-th step, also the sampling weights for the AUX data set $\boldsymbol{w}_{t+1}$. We created a video[1] visualizing the learning dynamics of Hopfield Boosting. For more information, we refer to Appendix C.4.
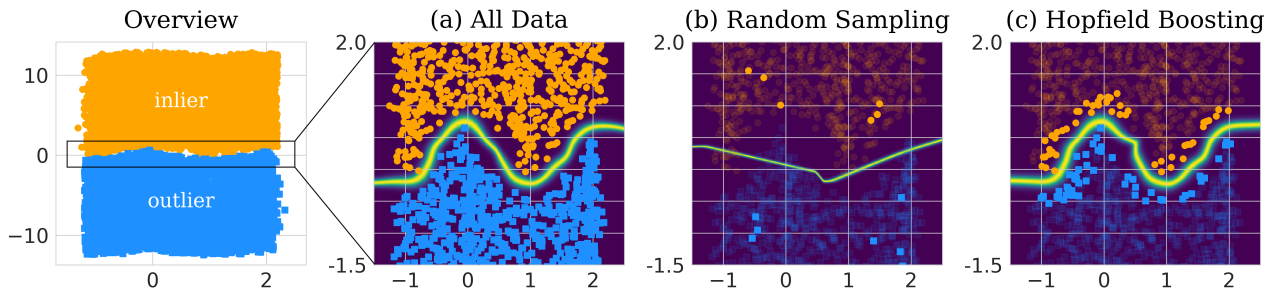
---

[1] https://youtu.be/H5tGdL-0fok

Figure 2: Synthetic example of the adaptive resampling mechanism. Hopfield Boosting forms a strong learner by sampling and combining a set of weak learners close to the decision boundary. The heatmap on the background shows $\exp(\beta\mathrm{E}_b(\boldsymbol{\xi}, \boldsymbol{X}, \boldsymbol{O}))$, where $\beta$ is 60. Only the sampled (i.e., highlighted) points serve as memories $\boldsymbol{X}$ and $\boldsymbol{O}$.



(a)          (b)

Figure 3: Depiction of the $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ energy landscape for separating inliers from AUX data / outliers on a 3D sphere with exemplary inlier (orange) and outlier (blue) points. (a) presents a front view, and (b) rotates the sphere by 90 degrees around the vertical axis. We set $\beta$ to 128. Appendix C.1 shows the sphere in more orientations.

**Inference.** At inference time, the OOD score $s(\boldsymbol{\xi})$ is

$$s(\boldsymbol{\xi}) = \mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) - \mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}). \quad (13)$$

For computing $s(\boldsymbol{\xi})$, Hopfield Boosting uses the *full* ID data set and a randomly sampled subset of the AUX data that has identical size. As we show in Appendix F.6, this step entails only a very moderate computational overhead in relation to a complete forward pass (e.g., an overhead of 7.5% for ResNet-18 on an NVIDIA Titan V GPU with 50,000 patterns stored in each of the memories $\boldsymbol{X}$ and $\boldsymbol{O}$).

We additionally experimented with using only $\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi})$ as a score, which also gives reasonable results. However, the approach in Equation (13) has turned out to be superior. Equation (13) uses information from both ID and AUX samples. This can, for example, be beneficial for handling query patterns $\boldsymbol{\xi}$ that are dissimilar from both the memory patterns in $\boldsymbol{X}$ as well as from the memory patterns in $\boldsymbol{O}$.

### 3.4. Comparison of Hopfield Boosting to HE and SHE

Zhang et al. (2023a) propose two post-hoc methods for OOD detection using MHE. They are called "Hopfield Energy" (HE) and "Simplified Hopfield Energy" (SHE). In contrast to Hopfield Boosting, HE and SHE do not make use of AUX data to get a better boundary between ID and OOD data. Rather, their methods evaluate the MHE on ID patterns only to determine whether a sample is ID or OOD. There are also more nuanced differences, for example, Hopfield Boosting uses a single Hopfield network for all ID classes, while HE and SHE employ an individual Hopfield network for each ID class. We provide a more in-depth discussion on the relation to HE and SHE in Appendix E.4.

## 4. Toy Examples

This section presents two toy examples illustrating the main intuitions behind Hopfield Boosting. For the sake of clarity, none of them considers the inlier classification task that would induce secondary processes, which would obscure the explanations. Formally, this means that we do not consider the first term on the right-hand side of Equation (11). For further toy examples, we refer to Appendix C.

The first example demonstrates how the weighting step in Hopfield Boosting allows good estimations of the decision boundary, even if Hopfield Boosting only samples a small number of weak learners (Figure 2). This is advantageous because the AUX data set contains a large number of data instances that are uninformative for the OOD detection task. For small, low dimensional data, one can always use all the data to compute $\mathrm{E}_b$ (Figure 2, a). For large problems (like in Ming et al., 2022), this strategy is difficult, and the naive solution of uniformly sampling N data points would also not work. This will yield many uninformative points (Figure 2, b). When using Hopfield Boosting and sampling N weak learners according to $\boldsymbol{w}_t$, the result better approximates the decision boundary of the full data (Figure 2, c).

Table 1: OOD detection performance on CIFAR-10. We compare results from Hopfield Boosting, DOS (Jiang et al., 2024), DOE (Wang et al., 2023b), DivOE (Zhu et al., 2023b), DAL (Wang et al., 2023a), MixOE (Zhang et al., 2023b), POEM (Ming et al., 2022), EBO-OE (Liu et al., 2020), and MSP-OE (Hendrycks et al., 2019b) on ResNet-18. ↓ indicates "lower is better" and ↑ "higher is better". All values in %. Standard deviations are estimated across five training runs.

| OOD Dataset | Metric | HB (ours) | DOS | DOE | DivOE | DAL | MixOE | POEM | EBO-OE | MSP-OE |
|---|---|---|---|---|---|---|---|---|---|---|
| SVHN | FPR95 ↓ | **0.23$^{\pm0.08}$** | 3.09$^{\pm0.75}$ | 1.97$^{\pm0.58}$ | 6.21$^{\pm0.84}$ | 1.25$^{\pm0.62}$ | 27.54$^{\pm2.46}$ | 1.48$^{\pm0.68}$ | 2.66$^{\pm0.91}$ | 4.31$^{\pm1.10}$ |
| | AUROC ↑ | **99.57$^{\pm0.06}$** | 99.15$^{\pm0.22}$ | 99.60$^{\pm0.13}$ | 98.53$^{\pm0.08}$ | **99.61$^{\pm0.15}$** | 95.37$^{\pm0.44}$ | 99.33$^{\pm0.15}$ | 99.15$^{\pm0.23}$ | 99.20$^{\pm0.15}$ |
| LSUN-Crop | FPR95 ↓ | 0.82$^{\pm0.17}$ | 3.66$^{\pm0.98}$ | 3.22$^{\pm0.45}$ | 1.88$^{\pm0.25}$ | 4.17$^{\pm0.27}$ | **0.14$^{\pm0.07}$** | 4.02$^{\pm0.91}$ | 6.82$^{\pm0.74}$ | 7.02$^{\pm1.14}$ |
| | AUROC ↑ | 99.40$^{\pm0.04}$ | 99.04$^{\pm0.20}$ | 99.30$^{\pm0.12}$ | 99.50$^{\pm0.02}$ | 99.13$^{\pm0.02}$ | **99.61$^{\pm0.11}$** | 98.89$^{\pm0.15}$ | 98.43$^{\pm0.10}$ | 98.83$^{\pm0.15}$ |
| LSUN-Resize | FPR95 ↓ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | 0.16$^{\pm0.17}$ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** |
| | AUROC ↑ | 99.98$^{\pm0.02}$ | 99.99$^{\pm0.01}$ | **100.00$^{\pm0.00}$** | 99.89$^{\pm0.05}$ | 99.92$^{\pm0.05}$ | 99.89$^{\pm0.06}$ | 99.88$^{\pm0.12}$ | 99.98$^{\pm0.02}$ | 99.96$^{\pm0.00}$ |
| Textures | FPR95 ↓ | **0.16$^{\pm0.02}$** | 1.28$^{\pm0.20}$ | 2.75$^{\pm0.57}$ | 1.20$^{\pm0.11}$ | 0.95$^{\pm0.13}$ | 4.68$^{\pm0.22}$ | 0.49$^{\pm0.04}$ | 1.11$^{\pm0.17}$ | 2.29$^{\pm0.16}$ |
| | AUROC ↑ | **99.84$^{\pm0.01}$** | 99.63$^{\pm0.04}$ | 99.35$^{\pm0.12}$ | 99.59$^{\pm0.02}$ | 99.74$^{\pm0.01}$ | 98.91$^{\pm0.07}$ | 99.72$^{\pm0.05}$ | 99.61$^{\pm0.02}$ | 99.57$^{\pm0.01}$ |
| iSUN | FPR95 ↓ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | 0.17$^{\pm0.12}$ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** |
| | AUROC ↑ | 99.97$^{\pm0.02}$ | 99.99$^{\pm0.01}$ | **100.00$^{\pm0.00}$** | 99.88$^{\pm0.05}$ | 99.93$^{\pm0.04}$ | 99.87$^{\pm0.05}$ | 99.87$^{\pm0.12}$ | 99.98$^{\pm0.01}$ | 99.96$^{\pm0.00}$ |
| Places 365 | FPR95 ↓ | **4.28$^{\pm0.23}$** | 12.26$^{\pm0.97}$ | 19.72$^{\pm2.39}$ | 13.70$^{\pm0.50}$ | 14.22$^{\pm0.51}$ | 16.30$^{\pm1.09}$ | 7.70$^{\pm0.68}$ | 11.77$^{\pm0.68}$ | 21.42$^{\pm0.88}$ |
| | AUROC ↑ | **98.51$^{\pm0.10}$** | 96.63$^{\pm0.43}$ | 95.06$^{\pm0.72}$ | 96.95$^{\pm0.09}$ | 96.77$^{\pm0.07}$ | 96.92$^{\pm0.22}$ | 97.56$^{\pm0.26}$ | 96.39$^{\pm0.30}$ | 95.91$^{\pm0.17}$ |
| Mean | FPR95 ↓ | **0.92** | 3.38 | 4.61 | 3.83 | 3.43 | 8.17 | 2.28 | 3.73 | 5.84 |
| | AUROC ↑ | **99.55** | 99.07 | 98.88 | 99.06 | 99.18 | 98.43 | 99.21 | 98.92 | 98.90 |

Table 2: OOD detection performance on CIFAR-100. We compare results from Hopfield Boosting, DOS (Jiang et al., 2024), DOE (Wang et al., 2023b), DivOE (Zhu et al., 2023b), DAL (Wang et al., 2023a), MixOE (Zhang et al., 2023b), POEM (Ming et al., 2022), EBO-OE (Liu et al., 2020), and MSP-OE (Hendrycks et al., 2019b) on ResNet-18. ↓ indicates "lower is better" and ↑ "higher is better". All values in %. Standard deviations are estimated across five training runs.

| OOD Dataset | Metric | HB (ours) | DOS | DOE | DivOE | DAL | MixOE | POEM | EBO-OE | MSP-OE |
|---|---|---|---|---|---|---|---|---|---|---|
| SVHN | FPR95 ↓ | 13.27$^{\pm5.46}$ | **9.84$^{\pm2.75}$** | 19.38$^{\pm4.60}$ | 28.77$^{\pm5.42}$ | 19.95$^{\pm2.34}$ | 41.54$^{\pm13.16}$ | 33.59$^{\pm4.12}$ | 36.33$^{\pm2.95}$ | 19.86$^{\pm6.90}$ |
| | AUROC ↑ | 97.07$^{\pm0.81}$ | **97.64$^{\pm0.39}$** | 95.72$^{\pm1.12}$ | 94.25$^{\pm0.98}$ | 95.69$^{\pm0.66}$ | 92.27$^{\pm2.71}$ | 94.06$^{\pm0.51}$ | 92.93$^{\pm0.72}$ | 95.74$^{\pm1.60}$ |
| LSUN-Crop | FPR95 ↓ | **12.68$^{\pm2.38}$** | 19.40$^{\pm2.45}$ | 28.23$^{\pm2.69}$ | 35.10$^{\pm4.23}$ | 24.24$^{\pm2.12}$ | 23.10$^{\pm7.39}$ | 15.72$^{\pm3.46}$ | 21.06$^{\pm3.12}$ | 32.88$^{\pm1.28}$ |
| | AUROC ↑ | **96.54$^{\pm0.65}$** | 96.42$^{\pm0.35}$ | 93.79$^{\pm0.88}$ | 92.45$^{\pm0.94}$ | 95.04$^{\pm0.43}$ | 96.11$^{\pm1.09}$ | **96.85$^{\pm0.60}$** | 95.79$^{\pm0.62}$ | 92.85$^{\pm0.33}$ |
| LSUN-Resize | FPR95 ↓ | **0.00$^{\pm0.00}$** | 0.01$^{\pm0.00}$ | 0.05$^{\pm0.04}$ | 0.01$^{\pm0.00}$ | **0.00$^{\pm0.00}$** | 10.27$^{\pm10.72}$ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | 0.03$^{\pm0.01}$ |
| | AUROC ↑ | 99.98$^{\pm0.01}$ | 99.96$^{\pm0.02}$ | 99.99$^{\pm0.01}$ | **99.99$^{\pm0.00}$** | 99.94$^{\pm0.02}$ | 97.99$^{\pm1.92}$ | 99.57$^{\pm0.09}$ | 99.57$^{\pm0.03}$ | 99.97$^{\pm0.00}$ |
| Textures | FPR95 ↓ | **2.35$^{\pm0.13}$** | 6.02$^{\pm0.52}$ | 19.42$^{\pm1.58}$ | 11.52$^{\pm0.49}$ | 5.22$^{\pm0.39}$ | 28.99$^{\pm6.79}$ | 2.89$^{\pm0.32}$ | 5.07$^{\pm0.54}$ | 10.34$^{\pm0.40}$ |
| | AUROC ↑ | **99.22$^{\pm0.02}$** | 98.33$^{\pm0.11}$ | 94.93$^{\pm0.48}$ | 97.02$^{\pm0.08}$ | 98.50$^{\pm0.16}$ | 94.24$^{\pm1.21}$ | 98.97$^{\pm0.08}$ | 98.15$^{\pm0.16}$ | 97.42$^{\pm0.08}$ |
| iSUN | FPR95 ↓ | **0.00$^{\pm0.00}$** | 0.03$^{\pm0.01}$ | 0.01$^{\pm0.00}$ | 0.06$^{\pm0.01}$ | 0.01$^{\pm0.02}$ | 14.40$^{\pm13.48}$ | **0.00$^{\pm0.00}$** | **0.00$^{\pm0.00}$** | 0.08$^{\pm0.02}$ |
| | AUROC ↑ | 99.98$^{\pm0.01}$ | 99.95$^{\pm0.02}$ | **99.99$^{\pm0.00}$** | 99.97$^{\pm0.00}$ | 99.93$^{\pm0.02}$ | 97.23$^{\pm2.59}$ | 99.59$^{\pm0.09}$ | 99.57$^{\pm0.03}$ | 99.96$^{\pm0.01}$ |
| Places 365 | FPR95 ↓ | 19.36$^{\pm1.02}$ | 32.13$^{\pm1.55}$ | 58.68$^{\pm4.15}$ | 44.20$^{\pm0.95}$ | 33.43$^{\pm1.11}$ | 47.01$^{\pm6.41}$ | **18.39$^{\pm0.68}$** | 26.68$^{\pm2.18}$ | 45.96$^{\pm0.85}$ |
| | AUROC ↑ | **95.85$^{\pm0.37}$** | 91.73$^{\pm0.39}$ | 83.47$^{\pm1.55}$ | 88.28$^{\pm0.26}$ | 91.10$^{\pm0.29}$ | 89.20$^{\pm1.86}$ | 95.03$^{\pm0.71}$ | 91.35$^{\pm0.70}$ | 87.77$^{\pm0.15}$ |
| Mean | FPR95 ↓ | **7.94** | 11.24 | 20.96 | 19.94 | 13.81 | 27.55 | 11.76 | 14.86 | 18.19 |
| | AUROC ↑ | **98.11** | 97.34 | 94.65 | 95.33 | 96.70 | 94.51 | 97.34 | 96.23 | 95.62 |

The second example depicts how inliers and outliers shape the energy surface (Figure 3). We generated patterns so that $X$ clusters around a pole and the outliers populate the remaining perimeter of the sphere. This is analogous to the idea that one has access to a large AUX data set, where some data points are more and some less informative for OOD detection (e.g., as conceptualized in Ming et al., 2022).

## 5. Experiments

### 5.1. Data & Setup

**CIFAR-10 & CIFAR-100.** Our training and evaluation proceeds as follows: We train Hopfield Boosting with ResNet-18 (He et al., 2016) on the CIFAR-10 and CIFAR-100 data sets (Krizhevsky, 2009), respectively. In these settings, we use ImageNet-RC (Chrabaszcz et al., 2017) (a low-resolution version of ImageNet) as the AUX data set. For testing the OOD detection performance, it is crucial to test how well the model can handle inputs from distributions it has never seen during training (i.e., that are different from the ID and the AUX data set): we use the data sets SVHN (Street View House Numbers) (Netzer et al., 2011), Textures (Cimpoi et al., 2014), iSUN (Xu et al., 2015), Places 365 (López-Cifuentes et al., 2020), and two versions of the LSUN data set (Yu et al., 2015) — one where the images are cropped, and one where they are resized to match the resolution of the CIFAR data sets (32x32 pixels). We compute the scores $s(\xi)$ as described in Equation (13) and then evaluate the discriminative power of $s(\xi)$ between CIFAR and the respective OOD data set using the FPR95 and the AUROC. We use a validation process with different OOD data for model selection. Specifically, we validate the model on MNIST (LeCun et al., 1998), and ImageNet-RC with

Table 3: OOD detection performance on ImageNet-1K. We compare results from Hopfield Boosting, DOS (Jiang et al., 2024), DOE (Wang et al., 2023b), DivOE (Zhu et al., 2023b), DAL (Wang et al., 2023a), MixOE (Zhang et al., 2023b), POEM (Ming et al., 2022), EBO-OE (Liu et al., 2020), and MSP-OE (Hendrycks et al., 2019b) on ResNet-50. ↓ indicates "lower is better" and ↑ "higher is better". All values in %. Standard deviations are estimated across five training runs.

| | | HB (ours) | DOS | DOE | DivOE | DAL | MixOE | POEM | EBO-OE | MSP-OE |
|---|---|---|---|---|---|---|---|---|---|---|
| Textures | FPR95 ↓ | $44.59^{\pm 1.05}$ | $40.29^{\pm 0.93}$ | $83.83^{\pm 7.19}$ | $42.80^{\pm 0.74}$ | $43.88^{\pm 0.66}$ | $41.05^{\pm 4.91}$ | $31.26^{\pm 0.67}$ | $\mathbf{29.67^{\pm 1.26}}$ | $48.38^{\pm 0.87}$ |
| | AUROC ↑ | $88.01^{\pm 0.57}$ | $89.88^{\pm 0.18}$ | $64.22^{\pm 9.25}$ | $88.18^{\pm 0.06}$ | $87.39^{\pm 0.15}$ | $88.51^{\pm 1.29}$ | $92.22^{\pm 0.90}$ | $\mathbf{92.40^{\pm 0.23}}$ | $86.25^{\pm 0.25}$ |
| SUN | FPR95 ↓ | $\mathbf{37.37^{\pm 1.84}}$ | $59.29^{\pm 0.96}$ | $83.73^{\pm 8.78}$ | $61.00^{\pm 0.57}$ | $65.31^{\pm 0.61}$ | $65.14^{\pm 2.53}$ | $57.46^{\pm 0.90}$ | $57.69^{\pm 1.61}$ | $66.01^{\pm 0.26}$ |
| | AUROC ↑ | $\mathbf{91.24^{\pm 0.52}}$ | $84.30^{\pm 0.21}$ | $72.95^{\pm 7.94}$ | $83.64^{\pm 0.30}$ | $81.47^{\pm 0.22}$ | $82.20^{\pm 0.72}$ | $85.38^{\pm 0.35}$ | $85.83^{\pm 0.60}$ | $81.45^{\pm 0.20}$ |
| Places 365 | FPR95 ↓ | $\mathbf{53.31^{\pm 2.05}}$ | $69.72^{\pm 1.01}$ | $86.30^{\pm 6.69}$ | $71.09^{\pm 0.60}$ | $74.46^{\pm 0.75}$ | $71.34^{\pm 1.49}$ | $68.87^{\pm 1.05}$ | $70.03^{\pm 1.83}$ | $74.58^{\pm 0.44}$ |
| | AUROC ↑ | $\mathbf{87.10^{\pm 0.52}}$ | $81.62^{\pm 0.22}$ | $70.37^{\pm 7.17}$ | $80.35^{\pm 0.33}$ | $78.72^{\pm 0.28}$ | $80.31^{\pm 0.42}$ | $81.79^{\pm 0.40}$ | $81.35^{\pm 0.63}$ | $78.89^{\pm 0.19}$ |
| iNaturalist | FPR95 ↓ | $\mathbf{11.11^{\pm 0.66}}$ | $49.55^{\pm 1.41}$ | $70.82^{\pm 13.89}$ | $30.51^{\pm 0.42}$ | $51.92^{\pm 0.74}$ | $47.28^{\pm 1.55}$ | $45.37^{\pm 1.79}$ | $49.02^{\pm 4.40}$ | $51.73^{\pm 1.35}$ |
| | AUROC ↑ | $\mathbf{97.65^{\pm 0.20}}$ | $90.49^{\pm 0.38}$ | $83.82^{\pm 5.75}$ | $93.81^{\pm 0.10}$ | $88.33^{\pm 0.21}$ | $90.19^{\pm 0.35}$ | $92.01^{\pm 0.33}$ | $91.44^{\pm 0.79}$ | $88.51^{\pm 0.30}$ |
| Mean | FPR95 ↓ | $\mathbf{36.60}$ | 54.71 | 81.17 | 51.35 | 58.90 | 56.20 | 50.74 | 51.60 | 60.17 |
| | AUROC ↑ | $\mathbf{91.00}$ | 86.57 | 72.84 | 86.49 | 83.98 | 85.30 | 87.85 | 87.75 | 83.78 |

Table 4: OOD detection performance on CIFAR-10. We compare Hopfield Boosting trained with weighted sampling and with random sampling on ResNet-18. ↓ indicates "lower is better" and ↑ indicates "higher is better". All values in %. Standard deviations are estimated across five training runs.

| | Weighted Sampling | | Random Sampling | |
|---|---|---|---|---|
| OOD Dataset | FPR95 ↓ | AUROC ↑ | FPR95 ↓ | AUROC ↑ |
| SVHN | $\mathbf{0.23^{\pm 0.08}}$ | $\mathbf{99.57^{\pm 0.06}}$ | $0.70^{\pm 0.13}$ | $\mathbf{99.55^{\pm 0.08}}$ |
| LSUN-Crop | $\mathbf{0.82^{\pm 0.20}}$ | $\mathbf{99.40^{\pm 0.05}}$ | $1.58^{\pm 0.31}$ | $99.24^{\pm 0.10}$ |
| LSUN-Resize | $\mathbf{0.00^{\pm 0.00}}$ | $\mathbf{99.98^{\pm 0.02}}$ | $0.00^{\pm 0.00}$ | $\mathbf{99.98^{\pm 0.01}}$ |
| Textures | $\mathbf{0.16^{\pm 0.02}}$ | $\mathbf{99.85^{\pm 0.01}}$ | $0.26^{\pm 0.06}$ | $99.81^{\pm 0.02}$ |
| iSUN | $\mathbf{0.00^{\pm 0.00}}$ | $99.97^{\pm 0.02}$ | $0.00^{\pm 0.00}$ | $\mathbf{99.99^{\pm 0.00}}$ |
| Places 365 | $\mathbf{4.28^{\pm 0.26}}$ | $\mathbf{98.51^{\pm 0.11}}$ | $6.20^{\pm 0.21}$ | $97.68^{\pm 0.21}$ |
| **Mean** | **0.92** | **99.55** | 1.46 | 99.38 |

different pre-processing than in training (resize to 32x32 pixels instead of crop to 32x32 pixels), as well as Gaussian and uniform noise.

**ImageNet-1K.** Following Wang et al. (2023b), we evaluate Hopfield Boosting on the large-scale benchmark: We use ImageNet-1K as ID data set and ImageNet-21K (Ridnik et al., 2021) as AUX data set. The OOD test data sets are Textures (Cimpoi et al., 2014), SUN (Xu et al., 2015), Places 365 (López-Cifuentes et al., 2020), and iNaturalist (Van Horn et al., 2018). In this setting, we fine-tune a pre-trained ResNet-50 using Hopfield Boosting.

**Baselines.** As mentioned earlier, previous works contain extensive experimental evidence that OE methods offer superior OOD detection performance compared to methods without OE (see e.g., Ming et al., 2022; Wang et al., 2023a). Our experiments in Appendix F.7 confirm this. Therefore, here we focus on a comprehensive comparison of Hopfield Boosting to eight OE methods: MSP-OE (Hendrycks et al., 2019b), EBO-OE (Liu et al., 2020), POEM (Ming et al., 2022), MixOE (Zhang et al., 2023b), DAL (Wang et al., 2023a), DivOE (Zhu et al., 2023b), DOE (Wang et al., 2023b) and DOS (Jiang et al., 2024).

**Training setup.** The network trains for 100 epochs (CIFAR-10/100) or 4 epochs (ImageNet-1K), respectively. In each epoch, the model processes the entire ID data set and a selection of AUX samples (sampled according to $\boldsymbol{w}_t$). We sample mini-batches of size 128 per data set, resulting in a combined batch size of 256. We evaluate the composite loss from Equation (11) for each resulting mini-batch and update the model accordingly. After an epoch, we update the sample weights, yielding $\boldsymbol{w}_{t+1}$. For efficiency reasons, we only compute the weights for 500,000 AUX data instances which we denote as $\Xi$. The weights of the remaining samples are set to 0. During the sample weight update, Hopfield Boosting does not compute gradients or update model parameters. The update of the sample weights $\boldsymbol{w}_{t+1}$ proceeds as follows: First, we fill the memories $\boldsymbol{X}$ and $\boldsymbol{O}$ with 50,000 ID samples and 50,000 AUX samples, respectively. Second, we use the obtained $\boldsymbol{X}$ and $\boldsymbol{O}$ to get the energy $\mathrm{E}_b(\Xi; \boldsymbol{X}, \boldsymbol{O})$ for each of the 500,000 AUX samples in $\Xi$ and compute $\boldsymbol{w}_{t+1}$ according to Equation (10). In the following epoch, Hopfield Boosting samples the mini-batches $\boldsymbol{O}_s^{\mathcal{D}}$ according to $\boldsymbol{w}_{t+1}$ with replacement. To allow the storage of even more patterns in the Hopfield memory during the weight update process, one could incorporate a vector similarity engine (e.g., Douze et al., 2024) into the process. This

would potentially allow a less noisy estimate of the sample weights. For the sake of simplicity, we did not opt to do this in our implementation of Hopfield Boosting. As we show in section 5.2, Hopfield Boosting achieves state-of-the-art OOD detection results and can scale to large datasets (ImageNet-1K) even without access to a similarity engine.

**Hyperparameters.** Like Yang et al. (2022), we use SGD with an initial learning rate of 0.1 and a weight decay of $5 \cdot 10^{-4}$. We decrease the learning rate during the training process with a cosine schedule (Loshchilov & Hutter, 2016). We apply these settings to all OOD detection methods that we test. For training Hopfield Boosting, we use a single value for $\beta$ throughout the training and evaluation process and for all OOD data sets. We tune the value of $\beta$ for each ID data set separately by selecting the value of $\beta$ from the set $\{2, 4, 8, 16, 32\}$ that performs best in the validation process. We select $\lambda$ — the weight for the OOD loss $\mathcal{L}_{\text{OOD}}$ — from $\{0.1, 0.25, 0.5, 1.0\}$. In our experiments, $\beta = 4$ and $\lambda = 0.5$ yields the best results for CIFAR-10 and CIFAR-100. For ImageNet-1K, we set $\beta = 32$ and $\lambda = 0.25$.

### 5.2. Results & Discussion

Table 1 summarizes the results for CIFAR-10. Hopfield Boosting achieves equal or better performance compared to the other methods regarding the FPR95 metric for all OOD data sets. It surpasses POEM (the previously best OOD detection approach with OE in our comparison), improving the mean FPR95 metric from 2.28 to 0.92. We observe that all methods tested perform worst on the Places 365 data set. To gain more insights regarding this behavior, we look at the data instances from the Places 365 data set that Hopfield Boosting trained on CIFAR-10 most confidently classifies as in-distribution (i.e., which receive the highest scores $s(\boldsymbol{\xi})$). Visual inspection shows that among those images, a large portion contains objects from semantic classes included in CIFAR-10 (e.g., airplanes, horses, automobiles). We refer to Appendix F.4 for more details.

We also investigate the influence of boosting (Table 4): The experiment shows that the sampling of weak learners contributes considerably to the performance of Hopfield Boosting. Although Hopfield Boosting demonstrates superior performance compared to POEM even without boosting, informative outlier sampling can beat this version on every dataset on the FPR95 metric. As in the other experiments, the results in LSUN-Resize and iSUN are nearly perfect. For additional ablations, we refer to Appendix F.2.

When subjecting Hopfield Boosting to data sets that were designed to show the weakness of OOD detection approaches (Appendix F.5), we identify instances where a substantial number of outliers are wrongly classified as inliers. Testing with EBO-OE yields comparable outcomes, indicating that

this phenomenon extends beyond Hopfield Boosting.

For CIFAR-100 (Table 2), Hopfield Boosting also surpasses POEM (the previously best method), improving the mean FPR95 from 11.76 to 7.95. On the SVHN data set, Hopfield Boosting improves the FPR95 metric the most, decreasing it from 33.59 to 13.27.

On ImageNet-1K (Table 3), Hopfield Boosting surpasses all methods in our comparison in terms of both mean FPR95 and mean AUROC. Compared to POEM (the previously best method) Hopfield Boosting improves the mean FPR95 from 50.74 to 36.60. This demonstrates that Hopfield Boosting scales surprisingly well to large-scale settings. This is surprising because, in this setting, we only store 100,000 samples in the Hopfield memory during inference (less than 1% of the samples of the ID and AUX data sets combined).

## 6. Conclusions & Outlook

We introduce Hopfield Boosting: an approach for OOD detection with OE. Hopfield Boosting uses an energy term to *boost* a classifier between inlier and outlier data by sampling weak learners that are close to the decision boundary. We illustrate how Hopfield Boosting shapes the energy surface to form a decision boundary. Additionally, we demonstrate how the boosting mechanism creates a sharper decision boundary than with random sampling. We compare Hopfield Boosting to eight modern OOD detection approaches using OE. Overall, Hopfield Boosting shows the best results.

Potential future work comprises: (a) There is a clear opportunity to improve the evaluation procedure for OOD detection. Specifically, it remains unclear how reliably the performance on specific data sets can indicate the general ability to detect OOD inputs. Therefore, we aim to develop an evaluation procedure that addresses these issues, allowing for a more nuanced assessment with a focus on relevant real-world scenarios. This could encompass not only an expansion of data but also the introduction of novel metrics and methodological improvements.

(b) A drawback of OE approaches is their reliance on AUX data. Although OE-based approaches improve the OOD detection capability in general, the selection of the AUX data is crucial, as it determines the characteristics of the decision boundary surrounding the inlier data. In practice, the use of AUX data can be prohibitive in domains where only a few or no outliers at all are available for training the model. Hence, a promising avenue is to explore strategies to more directly control the decision boundary by generating artificial outliers, such as through augmentation techniques.

## Acknowledgements

## References

Abbott, L. F. and Arian, Y. Storage capacity of generalized networks. *Physical review A*, 36(10):5091, 1987.

Ahn, S., Korattikara, A., and Welling, M. Bayesian posterior sampling via stochastic gradient Fisher scoring. In *Proceedings of the 29th International Coference on International Conference on Machine Learning*, pp. 1771–1778, Madison, WI, USA, 2012. Omnipress.

Auer, A., Gauch, M., Klotz, D., and Hochreiter, S. Conformal prediction for time series with modern Hopfield networks. *arXiv preprint arXiv:2303.12783*, 2023.

Baldi, P. and Venkatesh, S. S. Number of stable points for spin-glasses and neural networks of higher orders. *Physical Review Letters*, 58(9):913, 1987.

Bishop, C. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.

Bishop, C. M. Novelty detection and neural network validation. *IEE Proceedings-Vision, Image and Signal processing*, 141(4):217–222, 1994.

Breiman, L. Arcing the edge. Technical report, Citeseer, 1997.

Caputo, B. and Niemann, H. Storage capacity of kernel associative memories. In *Artificial Neural Networks—ICANN 2002: International Conference Madrid, Spain, August 28–30, 2002 Proceedings 12*, pp. 51–56. Springer, 2002.

Chen, H., Lee, Y., Sun, G., Lee, H., Maxwell, T., and Giles, C. L. High order correlation model for associative memory. In *AIP Conference Proceedings*, volume 151, pp. 86–99. American Institute of Physics, 1986.

Chen, J., Li, Y., Wu, X., Liang, Y., and Jha, S. Atom: Robustifying out-of-distribution detection using outlier mining. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III 21*, pp. 430–445. Springer, 2021.

Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.

Chrabaszcz, P., Loshchilov, I., and Hutter, F. A downsampled variant of imagenet as an alternative to the CIFAR datasets. *arXiv preprint arXiv:1707.08819*, 2017.

Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., , and Vedaldi, A. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.

Cortes, C. and Vapnik, V. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.

Demiriz, A., Bennett, K. P., and Shawe-Taylor, J. Linear programming boosting via column generation. *Machine Learning*, 46:225–254, 2002.

Djurisic, A., Bozanic, N., Ashok, A., and Liu, R. Extremely simple activation shaping for out-of-distribution detection. In *The Eleventh International Conference on Learning Representations*, 2023.

Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H. The faiss library. 2024.

Du, X., Wang, Z., Cai, M., and Li, Y. Vos: Learning what you don't know by virtual outlier synthesis. *arXiv preprint arXiv:2202.01197*, 2022.

Freund, Y. and Schapire, R. E. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational Learning Theory: Eurocolt '95*, pp. 23–37. Springer-Verlag, 1995.

Friedman, J., Hastie, T., and Tibshirani, R. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.

Fürst, A., Rumetshofer, E., Lehner, J., Tran, V. T., Tang, F., Ramsauer, H., Kreil, D., Kopp, M., Klambauer, G., Bitto, A., et al. CLOOB: Modern Hopfield networks with InfoLOOB outperform clip. *Advances in neural information processing systems*, 35:20450–20468, 2022.

Gardner, E. Multiconnected neural network models. *Journal of Physics A: Mathematical and General*, 20(11):3453, 1987.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.

Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Hkg4TI9xl.

Hendrycks, D., Basart, S., Mazeika, M., Zou, A., Kwon, J., Mostajabi, M., Steinhardt, J., and Song, D. Scaling out-of-distribution detection for real-world settings. *arXiv preprint arXiv:1911.11132*, 2019a.

Hendrycks, D., Mazeika, M., and Dietterich, T. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019b. URL https://openreview.net/forum?id=HyxCxhRcY7.

Hendrycks, D., Mazeika, M., Kadavath, S., and Song, D. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32, 2019c.

Hendrycks, D., Zou, A., Mazeika, M., Tang, L., Li, B., Song, D., and Steinhardt, J. Pixmix: Dreamlike pictures comprehensively improve safety measures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16783–16792, 2022.

Hopfield, J. J. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982.

Hopfield, J. J. Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences*, 81(10): 3088–3092, 1984. doi: 10.1073/pnas.81.10.3088.

Horn, D. and Usher, M. Capacities of multiconnected memory models. *Journal de Physique*, 49(3):389–395, 1988.

Huang, R., Geng, A., and Li, Y. On the importance of gradients for detecting distributional shifts in the wild. *Advances in Neural Information Processing Systems*, 34: 677–689, 2021.

Jiang, W., Cheng, H., Chen, M., Wang, C., and Wei, H. DOS: Diverse outlier sampling for out-of-distribution detection. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=iriEqxFB4y.

Krizhevsky, A. Learning multiple layers of features from tiny images. Master's thesis, Deptartment of Computer Science, University of Toronto, 2009.

Krotov, D. and Hopfield, J. J. Dense associative memory for pattern recognition. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, pp. 1172–1180. Curran Associates, Inc., 2016.

LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/file/abdeb6f575ac5c6676b747bca8d09cc2-Paper.pdf.

Liang, S., Li, Y., and Srikant, R. Enhancing the reliability of out-of-distribution image detection in neural networks. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=H1VGkIxRZ.

Liu, W., Wang, X., Owens, J., and Li, Y. Energy-based out-of-distribution detection. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 21464–21475. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/

f5496252609c43eb8a3d147ab9b9c006-Paper.pdf.

Liu, X., Lochman, Y., and Zach, C. Gen: Pushing the limits of softmax-based out-of-distribution detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23946–23955, 2023.

López-Cifuentes, A., Escudero-Vinolo, M., Bescós, J., and García-Martín, Á. Semantic-aware scene recognition. *Pattern Recognition*, 102:107256, 2020.

Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.

Lu, H., Gong, D., Wang, S., Xue, J., Yao, L., and Moore, K. Learning with mixture of prototypes for out-of-distribution detection. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=uNkKaD3MCs.

McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.

Ming, Y., Fan, Y., and Li, Y. POEM: Out-of-distribution detection with posterior sampling. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 15650–15665. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/ming22a.html.

Ming, Y., Sun, Y., Dia, O., and Li, Y. How to exploit hyperspherical embeddings for out-of-distribution detection? In *The Eleventh International Conference on Learning Representations*, 2023.

Moody, J. and Darken, C. J. Fast learning in networks of locally-tuned processing units. *Neural computation*, 1(2):281–294, 1989.

Morteza, P. and Li, Y. Provable guarantees for understanding out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7831–7840, 2022.

Müller, K.-R., Smola, A. J., Rätsch, G., Schölkopf, B., Kohlmorgen, J., and Vapnik, V. Predicting time series with support vector machines. In *International conference on artificial neural networks*, pp. 999–1004. Springer, 1997.

Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.

Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

Paischer, F., Adler, T., Patil, V., Bitto-Nemling, A., Holzleitner, M., Lehner, S., Eghbal-Zadeh, H., and Hochreiter, S. History compression via language models in reinforcement learning. In *International Conference on Machine Learning*, pp. 17156–17185. PMLR, 2022.

Park, G. Y., Kim, J., Kim, B., Lee, S. W., and Ye, J. C. Energy-based cross attention for Bayesian context update in text-to-image diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

Psaltis, D. and Park, C. H. Nonlinear discriminant functions and associative memories. In *AIP conference Proceedings*, volume 151, pp. 370–375. American Institute of Physics, 1986.

Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlović, M., Sandve, G. K., Greiff, V., Kreil, D., Kopp, M., Klambauer, G., Brandstetter, J., and Hochreiter, S. Hopfield networks is all you need. In *9th International Conference on Learning Representations (ICLR)*, 2021. URL https://openreview.net/forum?id=tL89RnzIiCd.

Rätsch, G., Onoda, T., and Müller, K.-R. Soft margins for AdaBoost. *Machine learning*, 42:287–320, 2001.

Ridnik, T., Ben-Baruch, E., Noy, A., and Zelnik-Manor, L. Imagenet-21k pretraining for the masses. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.

Roth, K., Pemula, L., Zepeda, J., Schölkopf, B., Brox, T., and Gehler, P. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14318–14328, 2022.

Ruff, L., Kauffmann, J. R., Vandermeulen, R. A., Montavon, G., Samek, W., Kloft, M., Dietterich, T. G., and Müller, K.-R. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 109(5):756–795, 2021.

Saleh, R. A. and Saleh, A. Statistical properties of the log-cosh loss function used in machine learning. *arXiv preprint arXiv:2208.04564*, 2022.

Sanchez-Fernandez, A., Rumetshofer, E., Hochreiter, S., and Klambauer, G. CLOOME: a new search engine unlocks bioimaging databases for queries with chemical structures. *bioRxiv*, pp. 2022–11, 2022.

Schäfl, B., Gruber, L., Bitto-Nemling, A., and Hochreiter, S. Hopular: Modern Hopfield networks for tabular data. *arXiv preprint arXiv:2206.00664*, 2022.

Sehwag, V., Chiang, M., and Mittal, P. Ssd: A unified framework for self-supervised outlier detection. *arXiv preprint arXiv:2103.12051*, 2021.

smeschke. Four Shapes. https://www.kaggle.com/datasets/smeschke/four-shapes/, 2018. URL https://www.kaggle.com/datasets/smeschke/four-shapes/.

Sun, Y. and Li, Y. Dice: Leveraging sparsification for out-of-distribution detection. In *European Conference on Computer Vision*, pp. 691–708. Springer, 2022.

Sun, Y., Guo, C., and Li, Y. React: Out-of-distribution detection with rectified activations. *Advances in Neural Information Processing Systems*, 34:144–157, 2021.

Sun, Y., Ming, Y., Zhu, X., and Li, Y. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning*, pp. 20827–20840. PMLR, 2022.

Tack, J., Mo, S., Jeong, J., and Shin, J. Csi: Novelty detection via contrastive learning on distributionally shifted instances. *Advances in neural information processing systems*, 33:11839–11852, 2020.

Tao, L., Du, X., Zhu, X., and Li, Y. Non-parametric outlier synthesis. *arXiv preprint arXiv:2303.02966*, 2023.

Teh, Y. W., Thiery, A. H., and Vollmer, S. J. Consistency and fluctuations for stochastic gradient Langevin dynamics. *J. Mach. Learn. Res.*, 17(1):193–225, 2016.

Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., and Belongie, S. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8769–8778, 2018.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc., 2017.

Wang, H., Li, Z., Feng, L., and Zhang, W. Vim: Out-of-distribution with virtual-logit matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4921–4930, 2022.

Wang, Q., Fang, Z., Zhang, Y., Liu, F., Li, Y., and Han, B. Learning to augment distributions for out-of-distribution detection. In *NeurIPS*, 2023a. URL https://openreview.net/forum?id=OtU6VvXJue.

Wang, Q., Ye, J., Liu, F., Dai, Q., Kalander, M., Liu, T., Hao, J., and Han, B. Out-of-distribution detection with implicit outlier transformation. *arXiv preprint arXiv:2303.05033*, 2023b.

Wang, T. and Isola, P. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pp. 9929–9939. PMLR, 2020.

Wei, H., Xie, R., Cheng, H., Feng, L., An, B., and Li, Y. Mitigating neural network overconfidence with logit normalization. In *International conference on machine learning*, pp. 23631–23644. PMLR, 2022a.

Wei, X.-S., Cui, Q., Yang, L., Wang, P., Liu, L., and Yang, J. Rpc: a large-scale and fine-grained retail product checkout dataset, 2022b. URL https://rpc-dataset.github.io/.

Welling, M. and Teh, Y. W. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning*, pp. 681–688, Madison, WI, USA, 2011. Omnipress.

Widrich, M., Schäfl, B., Pavlović, M., Ramsauer, H., Gruber, L., Holzleitner, M., Brandstetter, J., Sandve, G. K., Greiff, V., Hochreiter, S., and Klambauer, G. Modern Hopfield networks and attention for immune repertoire classification. *ArXiv*, 2007.13505, 2020.

Xu, K., Chen, R., Franchi, G., and Yao, A. Scaling for training time and post-hoc out-of-distribution detection enhancement. In *The Twelfth International Conference on Learning Representations*, 2024.

Xu, P., Ehinger, K. A., Zhang, Y., Finkelstein, A., Kulkarni, S. R., and Xiao, J. Turkergaze: Crowdsourcing saliency with webcam based eye tracking. *arXiv preprint arXiv:1504.06755*, 2015.

Xu, P., Chen, J., Zou, D., and Gu, Q. Global convergence of Langevin dynamics based algorithms for nonconvex optimization. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31, pp. 3122–3133. Curran Associates, Inc., 2018.

Yang, J., Wang, P., Zou, D., Zhou, Z., Ding, K., Peng, W., Wang, H., Chen, G., Li, B., Sun, Y., et al. Openood: Benchmarking generalized out-of-distribution detection. *Advances in Neural Information Processing Systems*, 35: 32598–32611, 2022.

Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., and Xiao, J. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.

Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6023–6032, 2019.

Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.

Zhang, J., Fu, Q., Chen, X., Du, L., Li, Z., Wang, G., Han, S., Zhang, D., et al. Out-of-distribution detection based on in-distribution data patterns memorization with modern Hopfield energy. In *The Eleventh International Conference on Learning Representations*, 2023a.

Zhang, J., Inkawhich, N., Linderman, R., Chen, Y., and Li, H. Mixture outlier exposure: Towards out-of-distribution detection in fine-grained environments. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 5531–5540, January 2023b.

Zhang, R., Li, C., Zhang, J., Chen, C., and Wilson, A. G. Cyclical stochastic gradient MCMC for Bayesian deep learning. In *International Conference on Learning Representations*, 2020.

Zheng, Y., Zhao, Y., Ren, M., Yan, H., Lu, X., Liu, J., and Li, J. Cartoon face recognition: A benchmark dataset. In *Proceedings of the 28th ACM International Conference on Multimedia*, pp. 2264–2272, 2020.

Zhu, F., Cheng, Z., Zhang, X.-Y., and Liu, C.-L. Openmix: Exploring outlier samples for misclassification detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12074–12083, 2023a.

Zhu, J., Geng, Y., Yao, J., Liu, T., Niu, G., Sugiyama, M., and Han, B. Diversified outlier exposure for out-of-distribution detection via informative extrapolation. *Advances in Neural Information Processing Systems*, 36, 2023b.

# A. Details on Continuous Modern Hopfield Networks

The following arguments are adopted from Fürst et al. (2022) and Ramsauer et al. (2021). Associative memory networks have been designed to store and retrieve samples. Hopfield networks are energy-based, binary associative memories, which were popularized as artificial neural network architectures in the 1980s (Hopfield, 1982; 1984). Their storage capacity can be considerably increased by polynomial terms in the energy function (Chen et al., 1986; Psaltis & Park, 1986; Baldi & Venkatesh, 1987; Gardner, 1987; Abbott & Arian, 1987; Horn & Usher, 1988; Caputo & Niemann, 2002; Krotov & Hopfield, 2016). In contrast to these binary memory networks, we use continuous associative memory networks with far higher storage capacity. These networks are continuous and differentiable, retrieve with a single update, and have exponential storage capacity (and are therefore scalable, i.e., able to tackle large problems; Ramsauer et al., 2021).

Formally, we denote a set of patterns $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\} \subset \mathbb{R}^d$ that are stacked as columns to the matrix $\boldsymbol{X} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N)$ and a state pattern (query) $\boldsymbol{\xi} \in \mathbb{R}^d$ that represents the current state. The largest norm of a stored pattern is $M = \max_i \|\boldsymbol{x}_i\|$. Then, the energy E of continuous Modern Hopfield Networks with state $\boldsymbol{\xi}$ is defined as (Ramsauer et al., 2021)

$$\mathrm{E} = -\beta^{-1} \log \left( \sum_{i=1}^{N} \exp(\beta \boldsymbol{x}_i^T \boldsymbol{\xi}) \right) + \frac{1}{2} \boldsymbol{\xi}^T \boldsymbol{\xi} + \mathrm{C}, \tag{14}$$

where $\mathrm{C} = \beta^{-1} \log N + \frac{1}{2} M^2$. For energy E and state $\boldsymbol{\xi}$, Ramsauer et al. (2021) proved that the update rule

$$\boldsymbol{\xi}^{\mathrm{new}} = \boldsymbol{X} \operatorname{softmax}(\beta \boldsymbol{X}^T \boldsymbol{\xi}) \tag{15}$$

converges globally to stationary points of the energy E and coincides with the attention mechanisms of Transformers (Vaswani et al., 2017; Ramsauer et al., 2021).

The *separation* $\Delta_i$ of a pattern $\boldsymbol{x}_i$ is its minimal dot product difference to any of the other patterns:

$$\Delta_i = \min_{j, j \neq i} \left( \boldsymbol{x}_i^T \boldsymbol{x}_i - \boldsymbol{x}_i^T \boldsymbol{x}_j \right). \tag{16}$$

A pattern is *well-separated* from the data if $\Delta_i$ is above a given threshold (specified in Ramsauer et al., 2021). If the patterns $\boldsymbol{x}_i$ are well-separated, the update rule Equation 15 converges to a fixed point close to a stored pattern. If some patterns are similar to one another and, therefore, not well-separated, the update rule converges to a fixed point close to the mean of the similar patterns.

The update rule of a Hopfield network thus identifies sample–sample relations between stored patterns. This enables similarity-based learning methods like nearest neighbor search (see Schäfl et al., 2022), which Hopfield Boosting leverages to detect samples outside the distribution of the training data.

# B. Notes on Langevin Sampling

Another method that is appropriate for earlier acquired models is to sample the posterior via the Stochastic Gradient Langevin Dynamics (SGLD) (Welling & Teh, 2011). This method is efficient since it iteratively learns from small mini-batches (Welling & Teh, 2011; Ahn et al., 2012). See basic work on Langevin dynamics (Welling & Teh, 2011; Ahn et al., 2012; Teh et al., 2016; Xu et al., 2018). A cyclical stepsize schedule for SGLD was very promising for uncertainty quantification (Zhang et al., 2020). Larger steps discover new modes, while smaller steps characterize each mode and perform the posterior sampling.

## C. Toy Examples

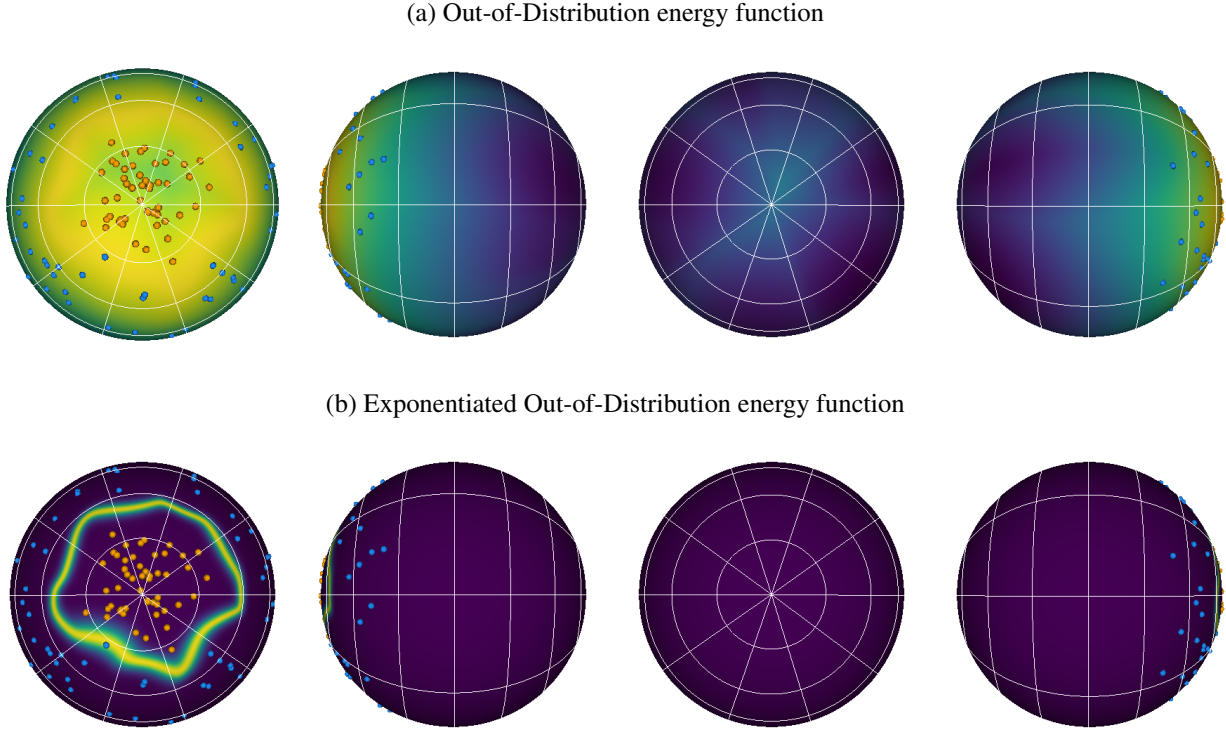### C.1. 3D Visualizations of $\mathrm{E}_b$ Loss in More Orientations

(a) Out-of-Distribution energy function



(b) Exponentiated Out-of-Distribution energy function



Figure 4: Depiction of the energy function $\mathrm{E}_b(\boldsymbol{\xi}, \boldsymbol{X}, \boldsymbol{O})$ in Fig. 3 in more orientations. (a) shows $\mathrm{E}_b(\boldsymbol{\xi}, \boldsymbol{X}, \boldsymbol{O})$ with exemplary inlier (orange) and outlier (blue) points; and (b) shows $\exp(\beta \mathrm{E}_b(\boldsymbol{\xi}, \boldsymbol{X}, \boldsymbol{O}))$. $\beta$ was set to 128. Both, (a) and (b), rotate the sphere by 0, 90, 180, and 270 degrees around the vertical axis.

### C.2. Dynamics of $\mathcal{L}_{\mathrm{OOD}}$ on Patterns in Euclidean Space

In this example, we applied our out-of-distribution loss $\mathcal{L}_{\mathrm{OOD}}$ on a simple binary classification problem. As we are working in Euclidean space and not on a sphere, we use a modified version of MHE, which uses the negative squared Euclidean distance instead of the dot-product-similarity. For the formal relation between Equation (17) and MHE, we refer to Appendix E.1:

$$\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X}) = -\beta^{-1} \log \left( \sum_{i=1}^{N} \exp(-\frac{\beta}{2} ||\boldsymbol{\xi} - \boldsymbol{x}_i||_2^2) \right) \tag{17}$$

Figure 5a shows the initial state of the patterns and the decision boundary $\exp(\beta E_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}))$. We store the samples of the two classes as stored patterns in $\boldsymbol{X}$ and $\boldsymbol{O}$, respectively, and compute $\mathcal{L}_{\mathrm{OOD}}$ for all samples. We then set the learning rate to 0.1 and perform gradient descent with $\mathcal{L}_{\mathrm{OOD}}$ on the data points. Figure 5b shows that after 25 steps, the distance between the data points and the decision boundary has increased, especially for samples that had previously been close to the decision boundary. After 100 steps, as shown in Figure 5d, the variability orthogonal to the decision boundary has almost completely vanished, while the variability parallel to the decision boundary is maintained.

(a) After 0 steps

(b) After 25 steps

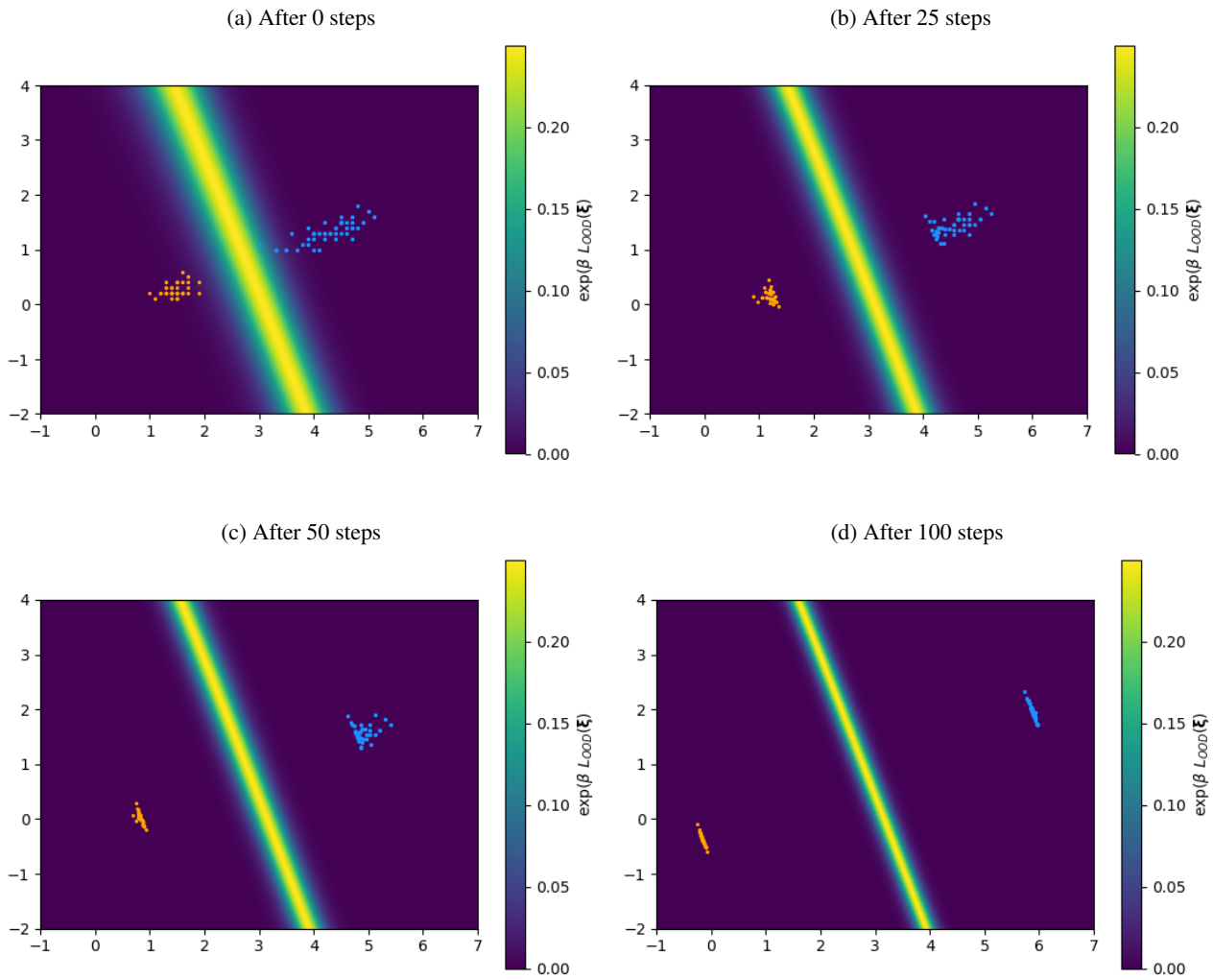(c) After 50 steps

(d) After 100 steps

Figure 5: $\mathcal{L}_{\text{OOD}}$ applied to exemplary data points on euclidean space. Gradient updates are applied to the data points directly. We observe that the variance orthogonal to the decision boundary shrinks while the variance parallel to the decision boundary does not change to this extent. $\beta$ is set to 2.
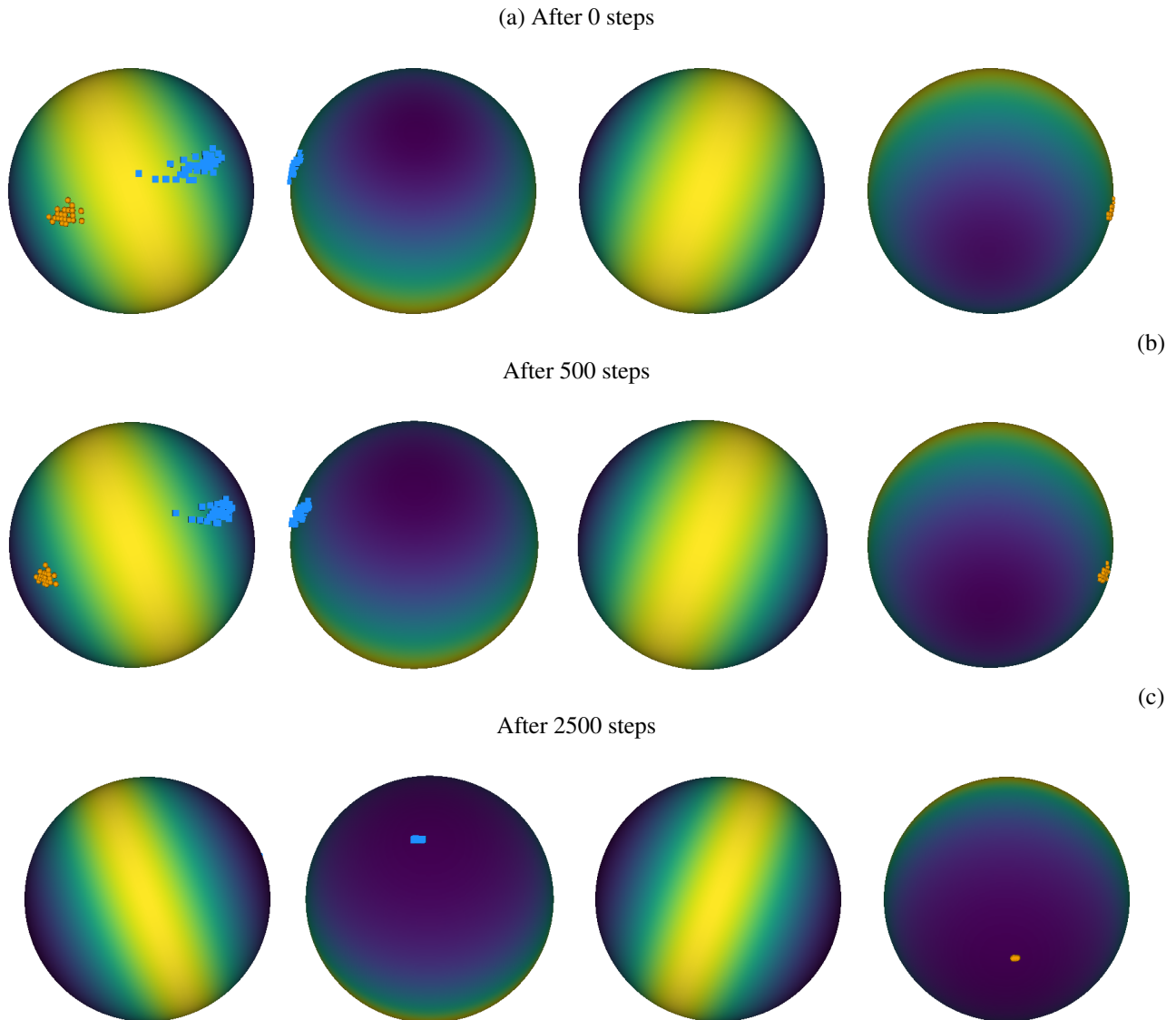
## C.3. Dynamics of $\mathcal{L}_{\text{OOD}}$ on Patterns on the Sphere

(a) After 0 steps



(b)

After 500 steps



(c)

After 2500 steps



Figure 6: $\mathcal{L}_{\text{OOD}}$ applied to exemplary data points on a sphere. Gradients are applied to the data points directly. We observe that the geometry of the space forces the patterns to opposing poles of the sphere.

## C.4. Learning Dynamics of Hopfield Boosting on Patterns on a Sphere - Video

The example video[2] demonstrates the learning dynamics of Hopfield Boosting on a 3-dimensional sphere. We randomly generate ID patterns $\boldsymbol{X}$ clustering around one of the sphere's poles and AUX patterns $\boldsymbol{O}$ on the remaining surface of the sphere. We then apply Hopfield Boosting on this data set. First, we sample the weak learners close to the decision boundary for both classes, $\boldsymbol{X}$ and $\boldsymbol{O}$. Then, we perform 2000 steps of gradient descent with $\mathcal{L}_{\text{OOD}}$ on the sampled weak learners. We apply the gradient updates to the patterns directly and do not propagate any gradients to an encoder. Every 50 gradient steps, we re-sample the weak learners. For this example, the initial learning rate is set to $0.02$ and increased after every gradient step by $0.1\%$.

---

[2]https://youtu.be/H5tGdL-0fok

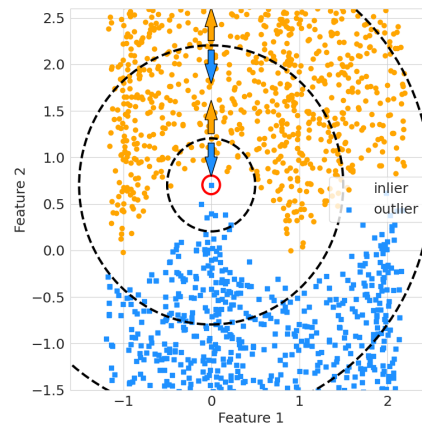## C.5. Location of Weak Learners near the Decision Boundary



Figure 7: A prototypical classifier (red circle) that is constructed with a sample close to the decision boundary. Classifiers like this one will only perform slightly better than random guessing (as indicated by the radial decision boundaries) and are, therefore, well-suited for weak learners.

## D. Notes on $\mathrm{E}_b$

### D.1. Probabilistic Interpretation of $\mathrm{E}_b$

We model the class-conditional densities of the in-distribution data and auxiliary data as mixtures of Gaussians with the patterns as the component means and tied, diagonal covariance matrices with $\beta^{-1}$ in the main diagonal.

$$p(\boldsymbol{\xi} \mid \mathrm{ID}) = \frac{1}{N} \sum_{i=1}^{N} \mathcal{N}\left(\boldsymbol{\xi}; \boldsymbol{x}_i, \beta^{-1} \boldsymbol{I}\right) \tag{18}$$

$$p(\boldsymbol{\xi} \mid \mathrm{AUX}) = \frac{1}{M} \sum_{i=1}^{M} \mathcal{N}\left(\boldsymbol{\xi}; \boldsymbol{o}_i, \beta^{-1} \boldsymbol{I}\right) \tag{19}$$

Further, we assume the distribution $p(\boldsymbol{\xi})$ as a mixture of $p(\boldsymbol{\xi} \mid \mathrm{ID})$ and $p(\boldsymbol{\xi} \mid \mathrm{AUX})$ with equal prior probabilities (mixture weights):

$$p(\boldsymbol{\xi}) = p(\mathrm{ID})\, p(\boldsymbol{\xi} \mid \mathrm{ID}) + p(\mathrm{AUX})\, p(\boldsymbol{\xi} \mid \mathrm{AUX}) \tag{20}$$

$$= \frac{1}{2}\, p(\boldsymbol{\xi} \mid \mathrm{ID}) + \frac{1}{2}\, p(\boldsymbol{\xi} \mid \mathrm{AUX}) \tag{21}$$

The probability of an unknown sample $\boldsymbol{\xi}$ being an AUX sample is given by

$$p(\mathrm{AUX} \mid \boldsymbol{\xi}) = \frac{p(\boldsymbol{\xi} \mid \mathrm{AUX})\, p(\mathrm{AUX})}{p(\boldsymbol{\xi})} \tag{22}$$

$$= \frac{p(\boldsymbol{\xi} \mid \mathrm{AUX})}{2\, p(\boldsymbol{\xi})} \tag{23}$$

$$= \frac{p(\boldsymbol{\xi} \mid \mathrm{AUX})}{p(\boldsymbol{\xi} \mid \mathrm{AUX}) + p(\boldsymbol{\xi} \mid \mathrm{ID})} \tag{24}$$

$$= \frac{1}{1 + \frac{p(\boldsymbol{\xi} \mid \mathrm{ID})}{p(\boldsymbol{\xi} \mid \mathrm{AUX})}} \tag{25}$$

$$= \frac{1}{1 + \exp(\log(p(\boldsymbol{\xi} \mid \mathrm{ID})) - \log(p(\boldsymbol{\xi} \mid \mathrm{AUX})))} \tag{26}$$

where in line (25) we have used that $p(\boldsymbol{\xi} \mid \mathrm{AUX}) > 0$ for all $\boldsymbol{\xi} \in \mathbb{R}^d$. The probability of $\boldsymbol{\xi}$ being an ID sample is given by

$$p(\mathrm{ID} \mid \boldsymbol{\xi}) = \frac{p(\boldsymbol{\xi} \mid \mathrm{ID})}{2\, p(\boldsymbol{\xi})} \tag{27}$$

$$= \frac{1}{1 + \exp(\log(p(\boldsymbol{\xi} \mid \mathrm{AUX})) - \log(p(\boldsymbol{\xi} \mid \mathrm{ID})))} \tag{28}$$

$$= 1 - p(\mathrm{AUX} \mid \boldsymbol{\xi}) \tag{29}$$

Consider the function

$$f_b(\boldsymbol{\xi}) = p(\mathrm{AUX} \mid \boldsymbol{\xi}) \cdot p(\mathrm{ID} \mid \boldsymbol{\xi}) \tag{30}$$

$$= \frac{p(\boldsymbol{\xi} \mid \mathrm{AUX}) \cdot p(\boldsymbol{\xi} \mid \mathrm{ID})}{4 p(\boldsymbol{\xi})^2} \tag{31}$$

By taking the $\log$ of Equation (31) we obtain the following. We use $\overset{C}{=}$ to denote equality up to an additive constant that does not depend on $\boldsymbol{\xi}$.

$$\beta^{-1} \log\left(f_b(\boldsymbol{\xi})\right) \overset{C}{=} -2\,\beta^{-1}\,\log\left(p(\boldsymbol{\xi})\right) \;+\; \beta^{-1}\,\log\left(p(\,\boldsymbol{\xi}\mid \text{ID}\,)\right) \;+\; \beta^{-1}\,\log\left(p(\,\boldsymbol{\xi}\mid \text{AUX}\,)\right) \tag{32}$$

Pre-multiplication by $\beta^{-1}$ is equivalent to a change of base of the log. The term $-\beta^{-1}\,\log(p(\boldsymbol{\xi}))$ is equivalent to the MHE (Ramsauer et al., 2021) (up to an additive constant) when assuming normalized patterns, i.e. $||\boldsymbol{x}_i||_2 = 1$ and $||\boldsymbol{o}_i||_2 = 1$, and an equal number of patterns $M = N$ in the two Gaussian mixtures $p(\,\boldsymbol{\xi}\mid \text{ID}\,)$ and $p(\,\boldsymbol{\xi}\mid \text{AUX}\,)$:

$$-\beta^{-1}\log(p(\boldsymbol{\xi})) = -\beta^{-1}\log\left(\frac{1}{2}p(\,\boldsymbol{\xi}\mid \text{ID}\,) \;+\; \frac{1}{2}p(\,\boldsymbol{\xi}\mid \text{AUX}\,)\right) \tag{33}$$

$$\overset{C}{=} -\beta^{-1}\log\left(p(\,\boldsymbol{\xi}\mid \text{ID}\,) \;+\; p(\,\boldsymbol{\xi}\mid \text{AUX}\,)\right) \tag{34}$$

$$= -\beta^{-1}\log\left(\frac{1}{N}\sum_{i=1}^{N}\mathcal{N}\left(\boldsymbol{\xi};\boldsymbol{x}_i,\beta^{-1}\boldsymbol{I}\right) \;+\; \frac{1}{N}\sum_{i=1}^{N}\mathcal{N}\left(\boldsymbol{\xi};\boldsymbol{o}_i,\beta^{-1}\boldsymbol{I}\right)\right) \tag{35}$$

$$\overset{C}{=} -\beta^{-1}\log\left(\sum_{i=1}^{N}\mathcal{N}\left(\boldsymbol{\xi};\boldsymbol{x}_i,\beta^{-1}\boldsymbol{I}\right) \;+\; \sum_{i=1}^{N}\mathcal{N}\left(\boldsymbol{\xi};\boldsymbol{o}_i,\beta^{-1}\boldsymbol{I}\right)\right) \tag{36}$$

$$\overset{C}{=} -\beta^{-1}\log\left(\sum_{i=1}^{N}\exp(-\frac{\beta}{2}||\boldsymbol{\xi}-\boldsymbol{x}_i||_2^2) \;+\; \sum_{i=1}^{N}\exp(-\frac{\beta}{2}||\boldsymbol{\xi}-\boldsymbol{o}_i||_2^2)\right) \tag{37}$$

$$\overset{C}{=} -\beta^{-1}\log\left(\sum_{i=1}^{N}\exp(\beta\boldsymbol{x}_i^T\boldsymbol{\xi}-\frac{\beta}{2}\boldsymbol{\xi}^T\boldsymbol{\xi}) \;+\; \sum_{i=1}^{N}\exp(\beta\boldsymbol{o}_i^T\boldsymbol{\xi}-\frac{\beta}{2}\boldsymbol{\xi}^T\boldsymbol{\xi})\right) \tag{38}$$

$$= -\beta^{-1}\log\left(\sum_{i=1}^{N}\exp(\beta\boldsymbol{x}_i^T\boldsymbol{\xi}) \;+\; \sum_{i=1}^{N}\exp(\beta\boldsymbol{o}_i^T\boldsymbol{\xi})\right) \;+\; \frac{1}{2}\,\boldsymbol{\xi}^T\boldsymbol{\xi} \tag{39}$$

$$\overset{C}{=} -\operatorname{lse}(\beta,(\boldsymbol{X}\,\|\,\boldsymbol{O})^T\boldsymbol{\xi}) \;+\; \frac{1}{2}\,\boldsymbol{\xi}^T\boldsymbol{\xi} \;+\; \beta^{-1}\log N + \frac{1}{2}M^2 \tag{40}$$

Analogously, $\beta^{-1}\,\log(p(\,\boldsymbol{\xi}\mid \text{ID}\,))$ and $\beta^{-1}\,\log(p(\,\boldsymbol{\xi}\mid \text{AUX}\,))$ also yield MHE terms. Therefore, $\mathrm{E}_b$ is equivalent to $\beta^{-1}\log(f_b(\boldsymbol{\xi}))$ under the assumption that $||\boldsymbol{x}_i||_2 = 1$ and $||\boldsymbol{o}_i||_2 = 1$ and $M = N$. The $\frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi}$ terms that are contained in the three MHEs cancel out.

$$\beta^{-1}\,\log\left(f_b(\boldsymbol{\xi})\right) \overset{C}{=} -2\operatorname{lse}(\beta,(\boldsymbol{X}\,\|\,\boldsymbol{O})^T\boldsymbol{\xi}) \;+\; \operatorname{lse}(\beta,\boldsymbol{X}^T\boldsymbol{\xi}) \;+\; \operatorname{lse}(\beta,\boldsymbol{O}^T\boldsymbol{\xi}) = \mathrm{E}_b(\boldsymbol{\xi};\boldsymbol{X},\boldsymbol{O}) \tag{41}$$

$f_b(\boldsymbol{\xi})$ can also be interpreted as the variance of a Bernoulli distribution with outcomes ID and AUX:

$$f_b(\boldsymbol{\xi}) = p(\,\text{AUX}\mid\boldsymbol{\xi}\,)\,p(\,\text{ID}\mid\boldsymbol{\xi}\,) = p(\,\text{ID}\mid\boldsymbol{\xi}\,)(1 - p(\,\text{ID}\mid\boldsymbol{\xi}\,)) = p(\,\text{AUX}\mid\boldsymbol{\xi}\,)(1 - p(\,\text{AUX}\mid\boldsymbol{\xi}\,)) \tag{42}$$

In other words, minimizing $\mathrm{E}_b$ means to drive a Bernoulli-distributed random variable with the outcomes ID and AUX towards minimum variance, i.e., $p(\,\text{ID}\mid\boldsymbol{\xi}\,)$ is driven towards 1 if $p(\,\text{ID}\mid\boldsymbol{\xi}\,) > 0.5$ and towards 0 if $p(\,\text{ID}\mid\boldsymbol{\xi}\,) < 0.5$. Conversely, the same is true for $p(\,\text{AUX}\mid\boldsymbol{\xi}\,)$.

From Equation (26), under the assumptions that $||\boldsymbol{x}_i||_2 = 1$ and $||\boldsymbol{o}_i||_2 = 1$ and $M = N$, the conditional probability $p(\,\text{AUX}\mid\boldsymbol{\xi}\,)$ can be computed as follows:

$$p(\,\text{AUX}\mid\boldsymbol{\xi}\,) = \sigma(\log(p(\,\boldsymbol{\xi}\mid \text{AUX}\,)) - \log(p(\,\boldsymbol{\xi}\mid \text{ID}\,))) \tag{43}$$

$$= \sigma(\beta\,(\operatorname{lse}(\beta,\boldsymbol{O}^T\boldsymbol{\xi}) - \operatorname{lse}(\beta,\boldsymbol{X}^T\boldsymbol{\xi}))) \tag{44}$$

where $\sigma$ denotes the logistic sigmoid function. Similarly, $p(\text{ID} \mid \boldsymbol{\xi})$ can be computed using

$$p(\text{ID} \mid \boldsymbol{\xi}) = \sigma(\beta \left(\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) - \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}))\right) \tag{45}$$

$$= 1 - p(\text{AUX} \mid \boldsymbol{\xi}) \tag{46}$$

## D.2. Alternative Formulations of $\text{E}_b$ and $f_b$

$\text{E}_b$ can be rewritten as follows.

$$\text{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) = -2\,\text{lse}(\beta, (\boldsymbol{X} \parallel \boldsymbol{O})^T \boldsymbol{\xi}) + \text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) + \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}) \tag{47}$$

$$= -2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}\left(\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) - \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}))\right)\right) - 2\beta^{-1}\,\log(2) \tag{48}$$

To prove this, we first show the following:

$$-\beta^{-1}\,\log\left(\exp(\beta\,\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi})) + \exp(\beta\,\text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}))\right) \tag{49}$$

$$= -\beta^{-1}\,\log\left(\exp\left(\beta\,\beta^{-1}\,\log\left(\sum_{i=1}^{N}\exp(\beta\boldsymbol{x}_i^T \boldsymbol{\xi})\right)\right) + \exp\left(\beta\,\beta^{-1}\,\log\left(\sum_{i=1}^{N}\exp(\beta\boldsymbol{o}_i^T \boldsymbol{\xi})\right)\right)\right) \tag{50}$$

$$= -\beta^{-1}\,\log\left(\sum_{i=1}^{N}\exp(\beta\boldsymbol{x}_i^T \boldsymbol{\xi}) + \sum_{i=1}^{N}\exp(\beta\boldsymbol{o}_i^T \boldsymbol{\xi})\right) \tag{51}$$

$$= -\text{lse}(\beta, (\boldsymbol{X} \parallel \boldsymbol{O})^T \boldsymbol{\xi}) \tag{52}$$

Let $\text{E}_{\boldsymbol{X}} = -\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi})$ and $\text{E}_{\boldsymbol{O}} = -\text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi})$.

$$\text{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) = -2\,\text{lse}(\beta, (\boldsymbol{X} \parallel \boldsymbol{O})^T \boldsymbol{\xi}) + \text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) + \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}) \tag{53}$$

$$= -2\beta^{-1}\,\log\left(\exp(-\beta\,\text{E}_{\boldsymbol{X}}) + \exp(-\beta\,\text{E}_{\boldsymbol{O}})\right) - \text{E}_{\boldsymbol{X}} - \text{E}_{\boldsymbol{O}} \tag{54}$$

$$= -2\beta^{-1}\,\log\left(\exp(-\frac{\beta}{2}\,\text{E}_{\boldsymbol{X}}) + \exp(-\beta\,\text{E}_{\boldsymbol{O}} + \frac{\beta}{2}\text{E}_{\boldsymbol{X}})\right) - \text{E}_{\boldsymbol{O}} \tag{55}$$

$$= -2\beta^{-1}\,\log\left(\exp(-\frac{\beta}{2}\,\text{E}_{\boldsymbol{X}} + \frac{\beta}{2}\,\text{E}_{\boldsymbol{O}}) + \exp(-\frac{\beta}{2}\,\text{E}_{\boldsymbol{O}} + \frac{\beta}{2}\,\text{E}_{\boldsymbol{X}})\right) \tag{56}$$

$$= -2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}(-\text{E}_{\boldsymbol{X}} + \text{E}_{\boldsymbol{O}})\right) - 2\beta^{-1}\,\log(2) \tag{57}$$

$$= -2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}\left(\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) - \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}))\right)\right) - 2\beta^{-1}\,\log(2) \tag{58}$$

$$= -2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}\left(\text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}) - \text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}))\right)\right) - 2\beta^{-1}\,\log(2) \tag{59}$$

By exponentiation of the above result we obtain

$$f_b(\boldsymbol{\xi}) \propto \exp(\beta\text{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})) = \frac{1}{4\cosh^2\left(\frac{\beta}{2}\left(\text{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) - \text{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}))\right)\right)} \tag{60}$$

The function $\log\cosh(x)$ is related to the negative log-likelihood of the hyperbolic secant distribution (see e.g. Saleh & Saleh, 2022). For values of $x$ close to 0, $\log\cosh$ can be approximated by $\frac{x^2}{2}$, and for values far from 0, the function behaves as $|x| - \log(2)$.
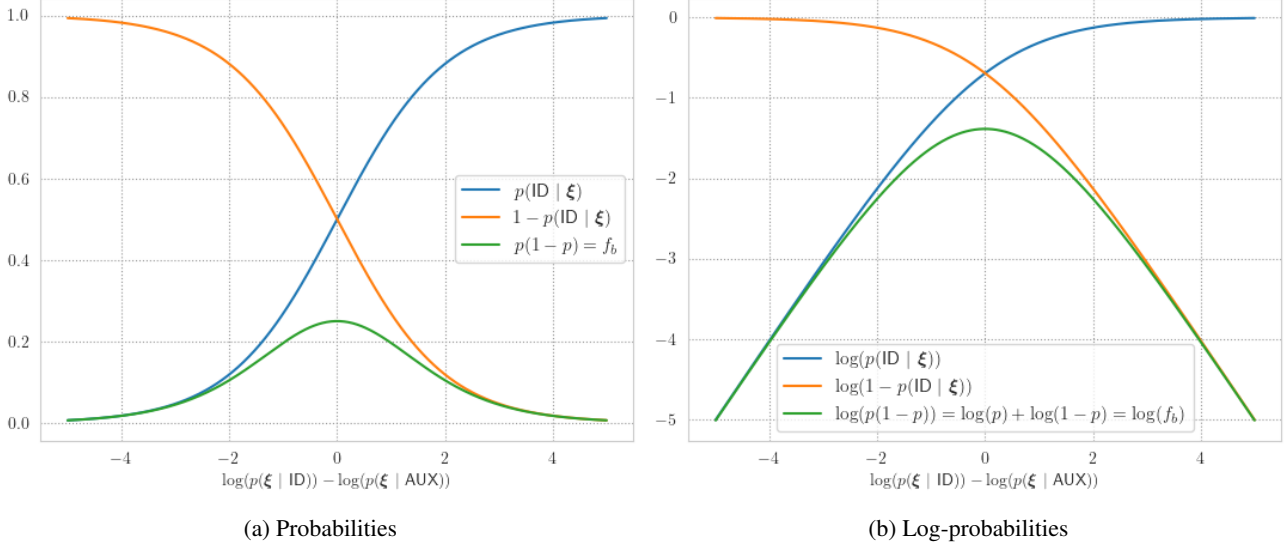
(a) Probabilities          (b) Log-probabilities

Figure 8: The product of two logistic sigmoids yields $f_b$ (a); the sum of two log-sigmoids yields $\log(f_b) = \mathrm{E}_b$ (b).

### D.3. Derivatives of $\mathrm{E}_b$

In this section, we investigate the derivatives of the energy function $\mathrm{E}_b$. The derivative of the lse is:

$$\nabla_{\boldsymbol{z}}\, \mathrm{lse}(\beta, \boldsymbol{z}) \;=\; \nabla_{\boldsymbol{z}}\, \beta^{-1}\, \log\left(\sum_{i=1}^{N} \exp(\beta z_i)\right) \;=\; \mathrm{softmax}(\beta\, \boldsymbol{z}) \tag{61}$$

Thus, the derivative of the MHE $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ w.r.t. $\boldsymbol{\xi}$ is:

$$\nabla_{\boldsymbol{\xi}}\, \mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X}) \;=\; \nabla_{\boldsymbol{\xi}} \left(-\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) + \frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + C\right) \;=\; -\boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}) + \boldsymbol{\xi} \tag{62}$$

The update rule of the MHN

$$\boldsymbol{\xi}^{t+1} \;=\; \boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}^t) \tag{63}$$

is derived via the concave-convex procedure. It coincides with the attention mechanisms of Transformers and has been proven to converge globally to stationary points of the energy $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ (Ramsauer et al., 2021). It can also be shown that the update rule emerges when performing gradient descent on $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ with step size $\eta = 1$ (Park et al., 2023):

$$\boldsymbol{\xi}^{t+1} \;=\; \boldsymbol{\xi}^t \;-\; \eta\, \nabla_{\boldsymbol{\xi}}\mathrm{E}(\boldsymbol{\xi}^t; \boldsymbol{X}) \tag{64}$$

$$\boldsymbol{\xi}^{t+1} \;=\; \boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}^t) \tag{65}$$

From Equation (62), we can see that the gradient of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ w.r.t. $\boldsymbol{\xi}$ is:

$$\nabla_{\boldsymbol{\xi}}\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) \;=\; \nabla_{\boldsymbol{\xi}}\left(-\,2\,\mathrm{lse}(\beta, (\boldsymbol{X} \,\|\, \boldsymbol{O})^T\boldsymbol{\xi}) \;+\; \mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) \;+\; \mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi})\right) \tag{66}$$

$$=\; -\,2\,(\boldsymbol{X} \,\|\, \boldsymbol{O})\,\mathrm{softmax}(\beta(\boldsymbol{X} \,\|\, \boldsymbol{O})^T\boldsymbol{\xi}) \;+\; \boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}) \;+\; \boldsymbol{O}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi}) \tag{67}$$

23

When $\boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi})$, $\boldsymbol{O}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi})$, $\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi})$ and $\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi})$ are available, one can efficiently compute $(\boldsymbol{X} \parallel \boldsymbol{O})\,\mathrm{softmax}(\beta(\boldsymbol{X} \parallel \boldsymbol{O})^T\boldsymbol{\xi})$ as follows:

$$(\boldsymbol{X} \parallel \boldsymbol{O})\,\mathrm{softmax}(\beta(\boldsymbol{X} \parallel \boldsymbol{O})^T\boldsymbol{\xi}) \;=\; \nabla_{\boldsymbol{\xi}}\,\mathrm{lse}(\beta, (\boldsymbol{X} \parallel \boldsymbol{O})^T\boldsymbol{\xi}) \tag{68}$$

$$= \nabla_{\boldsymbol{\xi}}\,\beta^{-1}\log\left(\exp(\beta\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi})) \;+\; \exp(\beta\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}))\right) \tag{69}$$

$$= \begin{pmatrix}\boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}) & \boldsymbol{O}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi})\end{pmatrix}\mathrm{softmax}\left(\beta\begin{pmatrix}\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi})\\ \mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi})\end{pmatrix}\right) \tag{70}$$

We can also compute the gradient of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ w.r.t. $\boldsymbol{\xi}$ via the $\log\cosh$-representation of $\mathrm{E}_b$ (see Equation (59)). The derivative of the $\log\cosh$ function is

$$\frac{\mathrm{d}}{\mathrm{d}x}\,\beta^{-1}\log\cosh(\beta x) \;=\; \tanh(\beta x) \tag{71}$$

Therefore, we can compute the gradient of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ as

$$\nabla_{\boldsymbol{\xi}}\,\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) = \nabla_{\boldsymbol{\xi}}\,-\,2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}\left(\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}) \,-\, \mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi})\right)\right) \tag{72}$$

$$= -\tanh\left(\frac{\beta}{2}(\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}) \,-\, \mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}))\right)\left(\boldsymbol{O}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi}) - \boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi})\right) \tag{73}$$

$$= -\tanh\left(\frac{\beta}{2}(\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) \,-\, \mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}))\right)\left(\boldsymbol{X}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi}) - \boldsymbol{O}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi})\right) \tag{74}$$

Next, we would like to compute the gradient of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ w.r.t. the memory matrices $\boldsymbol{X}$ and $\boldsymbol{O}$. For this, let us first look at the gradient of the MHE $\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X})$ w.r.t. a single stored pattern $\boldsymbol{x}_i$ (where $\boldsymbol{X}$ is the matrix of concatenated stored patterns $(\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N)$):

$$\nabla_{\boldsymbol{x}_i}\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X}) \;=\; -\,\boldsymbol{\xi}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi})_i \tag{75}$$

Thus, the gradient w.r.t. the full memory matrix $\boldsymbol{X}$ is

$$\nabla_{\boldsymbol{X}}\mathrm{E}(\boldsymbol{\xi}; \boldsymbol{X}) \;=\; -\boldsymbol{\xi}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi})^T \tag{76}$$

We can now also use the $\log\cosh$ formulation of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$ to compute the gradient of $\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$, w.r.t $\boldsymbol{X}$ and $\boldsymbol{O}$:

$$\nabla_{\boldsymbol{X}}\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) \;=\; \nabla_{\boldsymbol{X}}\,-\,2\beta^{-1}\,\log\cosh\left(\frac{\beta}{2}\left(\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) \,-\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}))\right)\right) \tag{77}$$

$$= -\,\tanh\left(\frac{\beta}{2}(\mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}) - \mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}))\right)\boldsymbol{\xi}\mathrm{softmax}(\beta\boldsymbol{X}^T\boldsymbol{\xi})^T \tag{78}$$

$$\tag{79}$$

Analogously, the gradient w.r.t $\boldsymbol{O}$ is

$$\nabla_{\boldsymbol{O}}\mathrm{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O}) \;=\; -\,\tanh\left(\frac{\beta}{2}(\mathrm{lse}(\beta, \boldsymbol{O}^T\boldsymbol{\xi}) - \mathrm{lse}(\beta, \boldsymbol{X}^T\boldsymbol{\xi}))\right)\boldsymbol{\xi}\mathrm{softmax}(\beta\boldsymbol{O}^T\boldsymbol{\xi})^T \tag{80}$$

# E. Notes on the Relationship between Hopfield Boosting and other methods

## E.1. Relation to Radial Basis Function Networks

This section shows the relation between radial basis function networks (RBF networks; Moody & Darken, 1989) and modern Hopfield energy (following Schäfl et al., 2022). Consider an RBF network with normalized linear weights:

$$\varphi(\boldsymbol{\xi}) = \sum_{i=1}^{N} \omega_i \exp(-\frac{\beta}{2} ||\boldsymbol{\xi} - \boldsymbol{\mu}_i||_2^2) \tag{81}$$

where $\beta$ denotes the inverse tied variance $\beta = \frac{1}{\sigma^2}$, and the $\omega_i$ are normalized using the $\mathrm{softmax}$ function:

$$\omega_i = \mathrm{softmax}(\beta \boldsymbol{a})_i = \frac{\exp(\beta a_i)}{\sum_{j=1}^{N} \exp(\beta a_j)} \tag{82}$$

An energy can be obtained by taking the negative log of $\varphi(\boldsymbol{\xi})$:

$$\begin{aligned} \mathrm{E}(\boldsymbol{\xi}) &= -\beta^{-1} \log\left(\varphi(\boldsymbol{\xi})\right) & (83) \\ &= -\beta^{-1} \log\left(\sum_{i=1}^{N} \omega_i \exp(-\frac{\beta}{2} ||\boldsymbol{\xi} - \boldsymbol{\mu}_i||_2^2))\right) & (84) \\ &= -\beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta(-\frac{1}{2}||\boldsymbol{\xi} - \boldsymbol{\mu}_i||_2^2 + \beta^{-1} \log \mathrm{softmax}(\beta \boldsymbol{a})_i))\right) & (85) \\ &= -\beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta(-\frac{1}{2}||\boldsymbol{\xi} - \boldsymbol{\mu}_i||_2^2 + a_i - \mathrm{lse}(\beta, \boldsymbol{a})))\right) & (86) \\ &= -\beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta(-\frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + \boldsymbol{\mu}_i^T\boldsymbol{\xi} - \frac{1}{2}\boldsymbol{\mu}_i^T\boldsymbol{\mu}_i + a_i))\right) + \mathrm{lse}(\beta, \boldsymbol{a}) & (87) \end{aligned}$$

Next, we define $a_i = \frac{1}{2}\boldsymbol{\mu}_i^T\boldsymbol{\mu}_i$

$$\mathrm{E}(\boldsymbol{\xi}) = -\beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta \boldsymbol{\mu}_i^T\boldsymbol{\xi})\right) + \frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + \mathrm{lse}(\beta, \boldsymbol{a}) \tag{88}$$

Finally, we use the fact that $\mathrm{lse}(\beta, \boldsymbol{a}) \le \max_i a_i + \beta^{-1} \log N$

$$\mathrm{E}(\boldsymbol{\xi}) = -\beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta \boldsymbol{\mu}_i^T\boldsymbol{\xi})\right) + \frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + \beta^{-1} \log N + \frac{1}{2}M^2 \tag{89}$$

where $M = \max_i ||\boldsymbol{\mu}_i||_2$

## E.2. Contrastive Representation Learning

A commonly used loss function in contrastive representation learning (e.g., Chen et al., 2020; He et al., 2020) is the InfoNCE loss (Oord et al., 2018):

$$\mathcal{L}_{\text{NCE}} = \mathop{\mathbb{E}}_{\substack{(x,y) \sim p_{\text{pos}} \\ \{x_i^-\}_{i=1}^M \sim p_{\text{data}}}} \left[ -\log \frac{e^{f(x)^T f(y)/\tau}}{e^{f(x)^T f(y)/\tau} + \sum_i e^{f(x_i^-)^T f(y)/\tau}} \right] \tag{90}$$

(Wang & Isola, 2020) show that $\mathcal{L}_{\text{NCE}}$ optimizes two objectives:

$$\mathcal{L}_{\text{NCE}} = \underbrace{\mathop{\mathbb{E}}_{(x,y) \sim p_{pos}} \left[ -f(x)^T f(y)/\tau \right]}_{\text{Alignment}} + \underbrace{\mathop{\mathbb{E}}_{\substack{(x,y) \sim p_{pos} \\ \{x_i^-\}_{i=1}^M \sim p_{data}}} \left[ \log \left( e^{f(x)^T f(y)/\tau} + \sum_i e^{f(x_i^-)^T f(x)/\tau} \right) \right]}_{\text{Uniformity}} \tag{91}$$

Alignment enforces that features from positive pairs are similar, while uniformity encourages a uniform distribution of the samples over the hypersphere.

In comparison, our proposed loss, $\mathcal{L}_{\text{OOD}}$, does not visibly enforce alignment between samples within the same class. Instead, we can observe that it promotes uniformity to the instances of the *foreign* class. Due to the constraints that are imposed by the geometry of the space the optimization is performed on, that is, $||f(x)|| = 1$ when the samples move on a hypersphere, the loss encourages the patterns in the ID data have maximum distance to the samples of the AUX data, i.e., they concentrate on opposing poles of the hypersphere. A demonstration of this mechanism can be found in Appendix C.2 and C.3

### E.3. Support Vector Machines

In the following, we will show the relation of Hopfield Boosting to support vector machines (SVMs; Cortes & Vapnik, 1995) with RBF kernel. We adopt and expand the arguments of Schäfl et al. (2022).

Assume we apply an SVM with RBF kernel to model the decision boundary between ID and AUX data. We train on the features $\boldsymbol{Z} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N, \boldsymbol{o}_1, \ldots, \boldsymbol{o}_M)$ and assume that the patterns are normalized, i.e., $||\boldsymbol{x}_i||_2 = 1$ and $||\boldsymbol{o}_i||_2 = 1$. We define the targets $(y_1, \ldots, y_{(N+M)})$ as 1 for ID and $-1$ for AUX data. The decision rule of the SVM equates to

$$\hat{B}(\boldsymbol{\xi}) = \begin{cases} \text{ID} & \text{if } s(\boldsymbol{\xi}) \geq 0 \\ \text{OOD} & \text{if } s(\boldsymbol{\xi}) < 0 \end{cases} \tag{92}$$

where

$$s(\boldsymbol{\xi}) = \sum_{i=1}^{N+M} \alpha_i y_i k(\boldsymbol{z}_i, \boldsymbol{\xi}) \tag{93}$$

$$k(\boldsymbol{z}_i, \boldsymbol{\xi}) = \exp\left( -\frac{\beta}{2} ||\boldsymbol{\xi} - \boldsymbol{z}_i||_2^2 \right) \tag{94}$$

We assume that there is at least one support vector for both ID and AUX data, i.e., there exists at least one index $i$ s.t. $\alpha_i y_i > 0$ and at least one index $j$ s.t. $\alpha_j y_j < 0$. We now split the samples $\boldsymbol{z}_i$ in $s(\boldsymbol{\xi})$ according to their label:

$$s(\boldsymbol{\xi}) = \sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi}) - \sum_{i=1}^{M} \alpha_{N+i} k(\boldsymbol{o}_i, \boldsymbol{\xi}) \tag{95}$$

We define an alternative score:

$$s_{\text{frac}}(\boldsymbol{\xi}) \; = \; \frac{\sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi})}{\sum_{i=1}^{M} \alpha_{N+i} k(\boldsymbol{o}_i, \boldsymbol{\xi})} \tag{96}$$

$$\tag{97}$$

Because we assumed there is at least one support vector for both ID and AUX data and as the $\alpha_i$ are constrained to be non-negative and because $k(\cdot, \cdot) > 0$, the numerator and denominator are strictly positive. We can, therefore, specify a new decision rule $\hat{B}_{\text{frac}}(\boldsymbol{\xi})$.

$$\hat{B}_{\text{frac}}(\boldsymbol{\xi}) \; = \; \begin{cases} \text{ID} & \text{if } s_{\text{frac}}(\boldsymbol{\xi}) \geq 1 \\ \text{OOD} & \text{if } s_{\text{frac}}(\boldsymbol{\xi}) < 1 \end{cases} \tag{98}$$

Although the functions $s(\boldsymbol{\xi})$ and $s_{\text{frac}}(\boldsymbol{\xi})$ are different, the decision rules $\hat{B}(\boldsymbol{\xi})$ and $\hat{B}_{\text{frac}}(\boldsymbol{\xi})$ are equivalent. Another possible pair of score and decision rule is the following:

$$s_{\log}(\boldsymbol{\xi}) \; = \; \beta^{-1} \log(s_{\text{frac}}(\boldsymbol{\xi})) \; = \; \beta^{-1} \log\left(\sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi})\right) \; - \; \beta^{-1} \log\left(\sum_{i=1}^{M} \alpha_{N+i} k(\boldsymbol{o}_i, \boldsymbol{\xi})\right) \tag{99}$$

$$\hat{B}_{\log}(\boldsymbol{\xi}) \; = \; \begin{cases} \text{ID} & \text{if } s_{\log}(\boldsymbol{\xi}) \geq 0 \\ \text{OOD} & \text{if } s_{\log}(\boldsymbol{\xi}) < 0 \end{cases} \tag{100}$$

Let us more closely examine the term $\beta^{-1} \log\left(\sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi})\right)$. We define $a_i = \beta^{-1} \log(\alpha_i)$.

$$\beta^{-1} \log\left(\sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi})\right) \; = \; \beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta a_i) \exp\left(-\frac{\beta}{2} \|\boldsymbol{\xi} - \boldsymbol{x}_i\|_2^2\right)\right) \tag{101}$$

$$= \; \beta^{-1} \log\left(\sum_{i=1}^{N} \exp(\beta a_i) \exp\left(-\frac{\beta}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + \beta\boldsymbol{x}_i^T\boldsymbol{\xi} - \frac{\beta}{2}\boldsymbol{x}_i^T\boldsymbol{x}_i\right)\right) \tag{102}$$

$$= \; \beta^{-1} \log\left(\sum_{i=1}^{N} \exp\left(-\frac{\beta}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} + \beta\boldsymbol{x}_i^T\boldsymbol{\xi} - \frac{\beta}{2}\boldsymbol{x}_i^T\boldsymbol{x}_i + \beta a_i\right)\right) \tag{103}$$

$$= \; \beta^{-1} \log\left(\sum_{i=1}^{N} \exp\left(\beta\boldsymbol{x}_i^T\boldsymbol{\xi} + \beta a_i\right)\right) \; - \; \frac{1}{2}\boldsymbol{\xi}^T\boldsymbol{\xi} \; - \; \frac{1}{2} \tag{104}$$

We now construct a memory $\boldsymbol{X}_H$ and query $\boldsymbol{\xi}_H$ such that we can compute (104) using the MHE (Equation (5)):

$$\boldsymbol{X}_H \; = \; \begin{pmatrix} \boldsymbol{x}_1 & \cdots & \boldsymbol{x}_N \\ a_1 & \cdots & a_N \end{pmatrix} \tag{105}$$

$$\boldsymbol{\xi}_H \; = \; \begin{pmatrix} \boldsymbol{\xi} \\ 1 \end{pmatrix} \tag{106}$$

We obtain

$$\mathrm{E}(\boldsymbol{\xi}_H; \boldsymbol{X}_H) \;=\; - \operatorname{lse}(\beta, \boldsymbol{X}_H^T \boldsymbol{\xi}_H) \;+\; \frac{1}{2}\boldsymbol{\xi}_H^T \boldsymbol{\xi}_H \;+\; C \tag{107}$$

$$= \; - \beta^{-1} \log \left( \sum_{i=1}^{N} \exp\left( \beta \boldsymbol{x}_i^T \boldsymbol{\xi} + 1\beta a_i \right) \right) \;+\; \frac{1}{2}\boldsymbol{\xi}^T \boldsymbol{\xi} \;+\; \frac{1}{2} \cdot 1^2 \;+\; C \tag{108}$$

$$= \; - \beta^{-1} \log \left( \sum_{i=1}^{N} \exp\left( \beta \boldsymbol{x}_i^T \boldsymbol{\xi} + \beta a_i \right) \right) \;+\; \frac{1}{2}\boldsymbol{\xi}^T \boldsymbol{\xi} \;+\; \frac{1}{2} \;+\; C \tag{109}$$

$$= \; - \beta^{-1} \log \left( \sum_{i=1}^{N} \alpha_i k(\boldsymbol{x}_i, \boldsymbol{\xi}) \right) \;+\; C \tag{110}$$

We construct $\boldsymbol{O}_H$ analogously to Equation (105) and thus can compute

$$s_{\log}(\boldsymbol{\xi}) \;=\; \mathrm{E}(\boldsymbol{\xi}_H; \boldsymbol{O}_H) \;-\; \mathrm{E}(\boldsymbol{\xi}_H; \boldsymbol{X}_H) \;=\; \operatorname{lse}(\beta, \boldsymbol{X}_H^T \boldsymbol{\xi}_H) \;-\; \operatorname{lse}(\beta, \boldsymbol{O}_H^T \boldsymbol{\xi}_H) \tag{111}$$

which is exactly the score Hopfield Boosting uses for determining whether a sample is OOD (Equation (13)). In contrast to SVMs, Hopfield Boosting uses a uniform weighting of the patterns in the memory when computing the score. However, Hopfield Boosting can emulate a weighting of the patterns by more frequently sampling patterns with high weights into the memory.

### E.4. HE and SHE

Zhang et al. (2023a) introduce two post-hoc methods for OOD detection using MHE, which are called "Hopfield Energy" (HE) and "Simplified Hopfield Energy" (SHE). Like Hopfield Boosting, HE and SHE both employ the MHE to determine whether a sample is ID or OOD. However, unlike Hopfield Boosting, HE and SHE offer no possibility to include AUX data in the training process to improve the OOD detection performance of their method. The rest of this section is structured as follows: First, we briefly introduce the methods HE and SHE, second, we formally analyze the two methods, and third, we formally relate them to Hopfield Boosting.

**Hopfield Energy (HE)** The method HE (Zhang et al., 2023a) computes the OOD score $s_{\mathrm{HE}}(\boldsymbol{\xi})$ as follows:

$$s_{\mathrm{HE}}(\boldsymbol{\xi}) \;=\; \operatorname{lse}(\beta, \boldsymbol{X}_c^T \boldsymbol{\xi}) \tag{112}$$

where $\boldsymbol{X}_c \in \mathbb{R}^{d \times N_c}$ denotes the memory $(\boldsymbol{x}_{c1}, \ldots, \boldsymbol{x}_{cN_c})$ containing $N_c$ encoded data instances of class $c$. HE uses the prediction of the ID classification head to determine which patterns to store in the Hopfield memory:

$$c \;=\; \operatorname*{argmax}_{y} p(\, y \mid \boldsymbol{\xi}^{\mathcal{D}} \,) \tag{113}$$

**Simplified Hopfield Energy (SHE)** The method SHE (Zhang et al., 2023a) employs a simplified score $s_{\mathrm{SHE}}(\boldsymbol{\xi})$:

$$s_{\mathrm{SHE}}(\boldsymbol{\xi}) \;=\; \boldsymbol{m}_c^T \boldsymbol{\xi} \tag{114}$$

where $\boldsymbol{m}_c \in \mathbb{R}^d$ denotes the mean of the patterns in memory $\boldsymbol{X}_c$:

$$\boldsymbol{m}_c \;=\; \frac{1}{N_c} \sum_{i=1}^{N_c} \boldsymbol{x}_{ci} \tag{115}$$

**Relation between HE and SHE**  In the following, we show a simple yet enlightening relation between the scores $s_{\text{HE}}$ and $s_{\text{SHE}}$. For mathematical convenience, we first slightly modify the score $s_{\text{HE}}$:

$$s_{\text{HE}}(\boldsymbol{\xi}) = \text{lse}(\beta, \boldsymbol{X}_c^T \boldsymbol{\xi}) - \beta^{-1} \log N_c \tag{116}$$

All data sets which were employed in the experiments of Zhang et al. (2023a) (CIFAR-10 and CIFAR-100) are class-balanced. Therefore, the additional term $\beta^{-1} \log N_c$ does not change the result of the OOD detection on those data sets, as it only amounts to the same constant offset for all classes.

The function

$$\text{lse}(\beta, \boldsymbol{z}) - \beta^{-1} \log N = \beta^{-1} \log \left( \frac{1}{N} \sum_{i=1}^{N} \exp(\beta z_i) \right) \tag{117}$$

converges to the mean function as $\beta \to 0$:

$$\lim_{\beta \to 0} \left( \text{lse}(\beta, \boldsymbol{z}) - \beta^{-1} \log N \right) = \frac{1}{N} \sum_{i=1}^{N} z_i \tag{118}$$

We now investigate the behavior of $s_{\text{HE}}$ in this limit:

$$\lim_{\beta \to 0} \left( \text{lse}(\beta, \boldsymbol{X}_c^T \boldsymbol{\xi}) - \beta^{-1} \log N \right) = \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{x}_{ci}^T \boldsymbol{\xi}) \tag{119}$$

$$= \left( \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_{ci} \right)^T \boldsymbol{\xi} \tag{120}$$

$$= \boldsymbol{m}_c^T \boldsymbol{\xi} \tag{121}$$

where

$$\boldsymbol{m}_c = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_{ci} \tag{122}$$

Therefore, we have shown that

$$\lim_{\beta \to 0} s_{\text{HE}}(\boldsymbol{\xi}) = s_{\text{SHE}}(\boldsymbol{\xi}) \tag{123}$$

**Relation of HE and SHE to Hopfield Boosting.**  In contrast to HE and SHE, Hopfield Boosting uses an AUX data set to learn a decision boundary between the ID and OOD regions during the training process. To do this, our work introduces a novel MHE-based energy function, $\text{E}_b(\boldsymbol{\xi}; \boldsymbol{X}, \boldsymbol{O})$, to determine how close a sample is to the learnt decision boundary. Hopfield Boosting uses this energy function to frequently sample weak learners into the Hopfield memory and for computing a novel Hopfield-based OOD loss $\mathcal{L}_{\text{OOD}}$. To the best our knowledge, we are the first to use MHE in this way to train a neural network.

The OOD detection score of Hopfield Boosting is

$$s(\boldsymbol{\xi}) \;=\; \mathrm{lse}(\beta, \boldsymbol{X}^T \boldsymbol{\xi}) \;-\; \mathrm{lse}(\beta, \boldsymbol{O}^T \boldsymbol{\xi}). \tag{124}$$

where $\boldsymbol{X} \in \mathbb{R}^{d \times N}$ contains the full encoded training set $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N)$ of all classes and $\boldsymbol{O} \in \mathbb{R}^{d \times M}$ contains AUX samples. While certainly similar to $s_{\mathrm{HE}}$, the Hopfield Boosting score $s$ differs from $s_{\mathrm{HE}}$ in three crucial aspects:

1. Hopfield Boosting uses AUX data samples in the OOD detection score in order to create a sharper decision boundary between the ID and OOD regions.

2. Hopfield Boosting normalizes the patterns in the memories $\boldsymbol{X}$ and $\boldsymbol{O}$ and the query $\boldsymbol{\xi}$ to unit length, while HE and SHE use unnormalized patterns to construct their memories $\boldsymbol{X}_c$ and their query pattern $\boldsymbol{\xi}$.

3. The score of Hopfield Boosting, $s(\boldsymbol{\xi})$, contains the full encoded training data set, while $s_{\mathrm{HE}}$ only contains the patterns of a single class. Therefore Hopfield Boosting computes the similarities of a query sample $\boldsymbol{\xi}$ to the entire ID data set. In Appendix F.6, we show that this process only incurs a moderate overhead of $7.5\%$ compared to the forward pass of the ResNet-18.

The selection of the score function $s(\boldsymbol{\xi})$ is only a small aspect of Hopfield Boosting. Hopfield Boosting additionally samples informative AUX data close to the decision boundary, optimizes an MHE-based loss function, and thereby learns a sharp decision boundary between ID and OOD regions. Those three aspects are novel contributions of Hopfield Boosting. In contrast, the work of Zhang et al. (2023a) solely focuses on the selection of a suitable Hopfield-based OOD detection score for post-hoc OOD detection.

Table 5: Comparison between HE, SHE and our version. ↓ indicates "lower is better" and ↑ indicates "higher is better".

| OOD Dataset | Ours FPR95 ↓ | Ours AUROC ↑ | HE FPR95 ↓ | HE AUROC ↑ | SHE FPR95 ↓ | SHE AUROC ↑ |
|---|---|---|---|---|---|---|
| SVHN | 36.79 | **93.18** | 35.81 | 92.35 | **35.07** | 92.81 |
| LSUN-Crop | **13.10** | **97.25** | 17.74 | 95.96 | 18.19 | 96.10 |
| LSUN-Resize | **16.65** | **96.84** | 20.69 | 95.87 | 21.66 | 95.85 |
| Textures | **44.54** | **89.38** | 46.29 | 86.67 | 46.19 | 87.44 |
| iSUN | **19.20** | **96.08** | 22.52 | 95.08 | 23.25 | 95.06 |
| Places 365 | **39.02** | **90.63** | 41.56 | 88.41 | 42.57 | 88.38 |
| **Mean** | **28.21** | **93.89** | 30.77 | 92.39 | 31.66 | 92.60 |

Table 6: Comparison of OOD detection performance on CIFAR-10 of Hopfield Boosting on different encoders. ↓ indicates "lower is better" and ↑ indicates "higher is better". Standard deviations are estimated across five independent training runs.

| OOD Dataset | ResNet-18 FPR95 ↓ | ResNet-18 AUROC ↑ | ResNet-34 FPR95 ↓ | ResNet-34 AUROC ↑ | ResNet-50 FPR95 ↓ | ResNet-50 AUROC ↑ | Densenet-100 FPR95 ↓ | Densenet-100 AUROC ↑ |
|---|---|---|---|---|---|---|---|---|
| SVHN | $0.23^{\pm 0.08}$ | $99.57^{\pm 0.06}$ | $0.33^{\pm 0.25}$ | $99.63^{\pm 0.07}$ | $\mathbf{0.19^{\pm 0.09}}$ | $\mathbf{99.64^{\pm 0.11}}$ | $2.11^{\pm 2.76}$ | $99.31^{\pm 0.35}$ |
| LSUN-Crop | $0.82^{\pm 0.20}$ | $99.40^{\pm 0.05}$ | $0.65^{\pm 0.14}$ | $99.54^{\pm 0.07}$ | $0.69^{\pm 0.15}$ | $99.47^{\pm 0.09}$ | $\mathbf{0.40^{\pm 0.23}}$ | $\mathbf{99.52^{\pm 0.09}}$ |
| LSUN-Resize | $\mathbf{0.00^{\pm 0.00}}$ | $99.98^{\pm 0.02}$ | $\mathbf{0.00^{\pm 0.00}}$ | $99.89^{\pm 0.04}$ | $\mathbf{0.00^{\pm 0.00}}$ | $99.93^{\pm 0.10}$ | $\mathbf{0.00^{\pm 0.00}}$ | $\mathbf{100.0^{\pm 0.00}}$ |
| Textures | $0.16^{\pm 0.02}$ | $99.85^{\pm 0.01}$ | $0.15^{\pm 0.07}$ | $99.89^{\pm 0.04}$ | $0.16^{\pm 0.07}$ | $99.83^{\pm 0.01}$ | $\mathbf{0.08^{\pm 0.03}}$ | $\mathbf{99.88^{\pm 0.01}}$ |
| iSUN | $\mathbf{0.00^{\pm 0.00}}$ | $99.97^{\pm 0.02}$ | $\mathbf{0.00^{\pm 0.00}}$ | $99.98^{\pm 0.02}$ | $\mathbf{0.00^{\pm 0.00}}$ | $99.98^{\pm 0.02}$ | $\mathbf{0.00^{\pm 0.00}}$ | $\mathbf{99.99^{\pm 0.01}}$ |
| Places 365 | $4.28^{\pm 0.26}$ | $98.51^{\pm 0.11}$ | $4.13^{\pm 0.54}$ | $98.46^{\pm 0.22}$ | $4.75^{\pm 0.45}$ | $98.71^{\pm 0.05}$ | $\mathbf{2.56^{\pm 0.20}}$ | $\mathbf{99.26^{\pm 0.03}}$ |
| **Mean** | 0.92 | 99.55 | 0.88 | 99.57 | 0.97 | 99.59 | **0.86** | **99.66** |

# F. Additional Experiments & Experimental Details

### F.1. Comparison HE/SHE

Since Hopfield Boosting shares similarities with the MHE-based methods HE and SHE (Zhang et al., 2023a), we also looked at the approach as used for their methods. We use the same ResNet-18 as a backbone network as we used in the experiments for Hopfield Boosting, but train it on CIFAR-10 without OE. We modify the approach of Zhang et al. (2023a) to not only use the penultimate layer, but perform a search over all layer activation combinations of the backbone for the best-performing combination. We also do not use the classifier to separate by class. From the search, we see that the concatenated activations of layers 3 and 5 give the best performance on average, so we use this setting. We experience a quite noticeable drop in performance compared to their results (Table 5). Since the computation of the MHE is the same, we assume the reason for the performance drop is the different training of the ResNet-18 backbone network, where (Zhang et al., 2023a) used strong augmentations.

### F.2. Ablations

We investigate the impact of different encoder backbone architectures on OOD detection performance with Hopfield Boosting. The baseline uses a ResNet-18 as the encoder architecture. For the ablation, the following architectures are used as a comparison: ResNet-34, ResNet-50, and Densenet-100. It can be observed, that the larger architectures lead to a slight increase in OOD performance (Table 6). We also see that a change in architecture from ResNet to Densenet leads to a different OOD behavior: The result on the Places365 data set is greatly improved, while the performance on SVHN is noticeably worse than on the ResNet architectures. The FPR95 of Densenet on SVHN also shows a high variance, which is due to one of the five independent training runs performing very badly at detecting SVHN samples as OOD: The worst run scores an FPR95 5.59, while the best run achieves an FPR95 of 0.24.
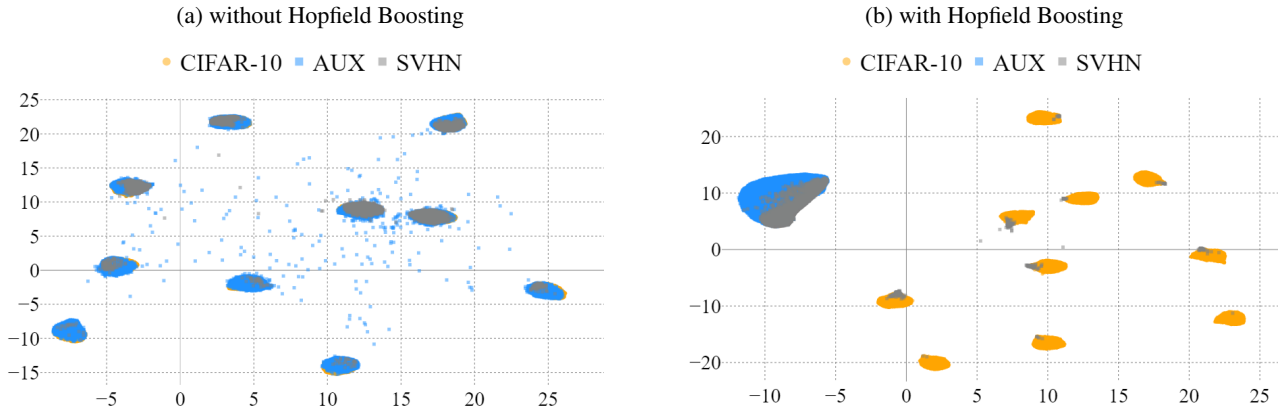
Figure 9: UMAP embeddings of ID (CIFAR-10) and OOD (AUX and SVHN) data based on our model trained without (a) and with Hopfield Boosting (b). Clearly, without Hopfield Boosting, the embedded OOD data points tend to overlap with the ID data points, making it impossible to distinguish between ID and OOD. On the other hand, Hopfield Boosting shows a clear separation of ID and OOD data in the embedding.

### F.3. Effect on Learned Representation

In order to analyze the impact of Hopfield Boosting on learned representations, we utilize the output of our model's embedding layer (see 5.1) as the input for a manifold learning-based visualization. Uniform Manifold Approximation and Projection (UMAP) McInnes et al. (2018) is a non-linear dimensionality reduction technique known for its ability to preserve both global and local structure in high-dimensional data.

First, we train two models – with and without Hopfield Boosting– and extract the embeddings of both ID and OOD data sets from them. This results in a 512-dimensional vector representation for each data point, which we further reduce to two dimensions with UMAP. The training data for UMAP always corresponds to the training data of the respective method. That is, the model trained without Hopfield Boosting is solely trained on CIFAR-10 data, and the model trained with Hopfield Boosting is presented with CIFAR-10 and AUX data during training, respectively. We then compare the learned representations concerning ID and OOD data.

Figure 9 shows the UMAP embeddings of ID (CIFAR-10) and OOD (AUX and SVHN) data based on our model trained without (a) and with Hopfield Boosting (b). Without Hopfield Boosting, OOD data points typically overlap with ID data points, with just a few exceptions, making it difficult to differentiate between them. Conversely, Hopfield Boosting allows to distinctly separate ID and OOD data in the embedding.

### F.4. OOD Examples from the Places 365 Data Set with High Semantic Similarity to CIFAR-10

We observe that Hopfield Boosting and all competing methods struggle with correctly classifying the samples from the Places 365 data set as OOD the most. Table 1 shows that for Hopfield Boosting, the FPR95 for the Places 365 data set with CIFAR-10 as the ID data set is at 4.28. The second worst FPR95 for Hopfield Boosting was measured on the LSUN-Crop data set at 0.82.

We inspect the 100 images from Places 365 that perform worst (i.e., that achieve the highest score $s(\boldsymbol{\xi})$) on a model trained with Hopfield Boosting on the CIFAR-10 data set as the in-distribution data set. Figure 10 shows that within those 100 images, the Places 365 data set contains a non-negligible amount of data instances that show objects from semantic classes contained in CIFAR-10 (e.g., horses, automobiles, dogs, trucks, and airplanes). We argue that data instances that clearly show objects of semantic classes contained in CIFAR-10 should be considered as in-distribution, which Hopfield Boosting correctly recognizes. Therefore, a certain amount of error can be anticipated on the Places 365 data set for all OOD detection methods. We leave a closer evaluation of the amount of the anticipated error up to future work.

For comparison, Figure 11 shows the 100 images from Places 365 with the lowest score $s(\boldsymbol{\xi})$, as evaluated by a model trained

with Hopfield Boosting on CIFAR-10. There are no objects visible that have clear semantic overlap with the CIFAR-10 classes.
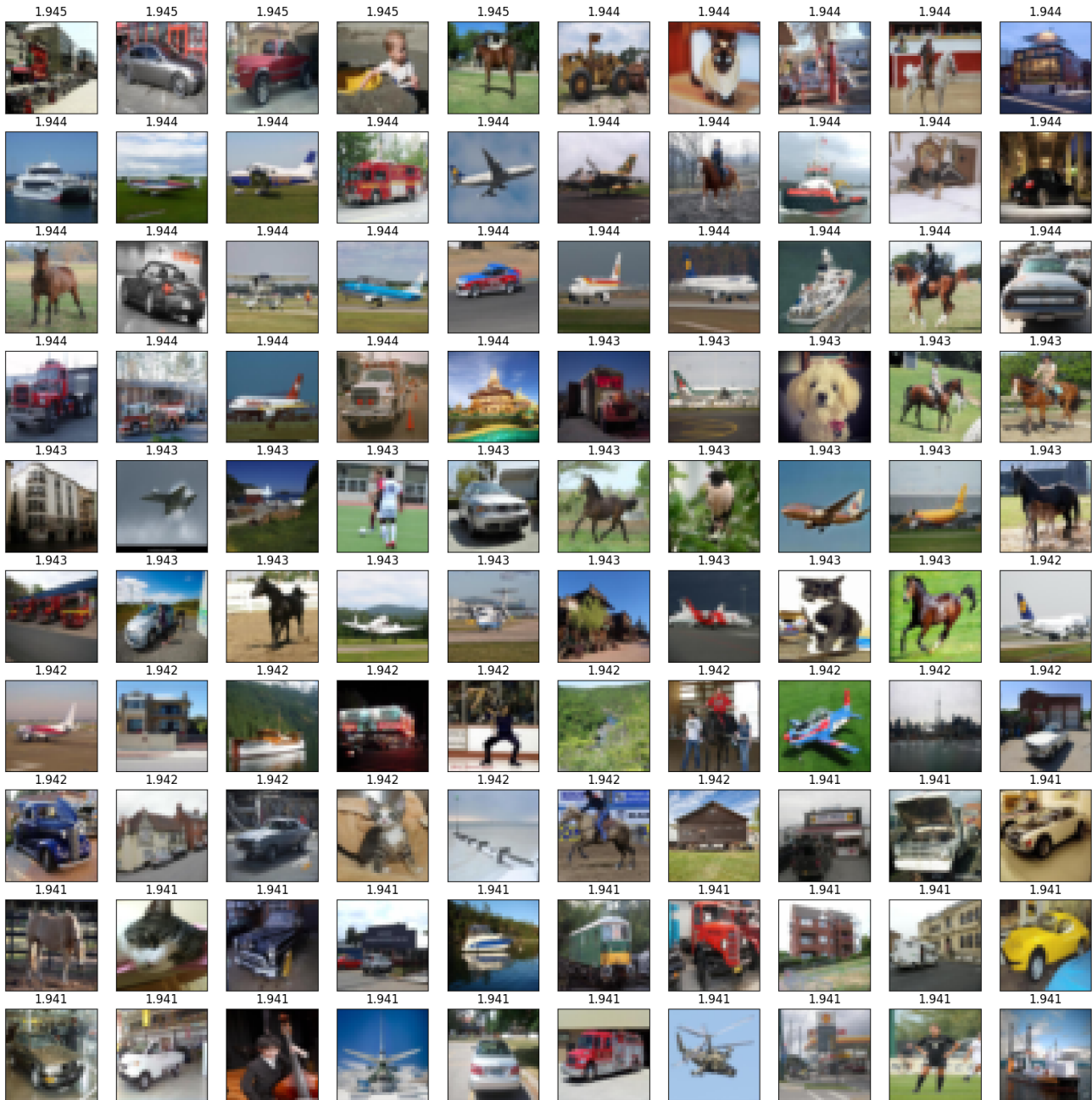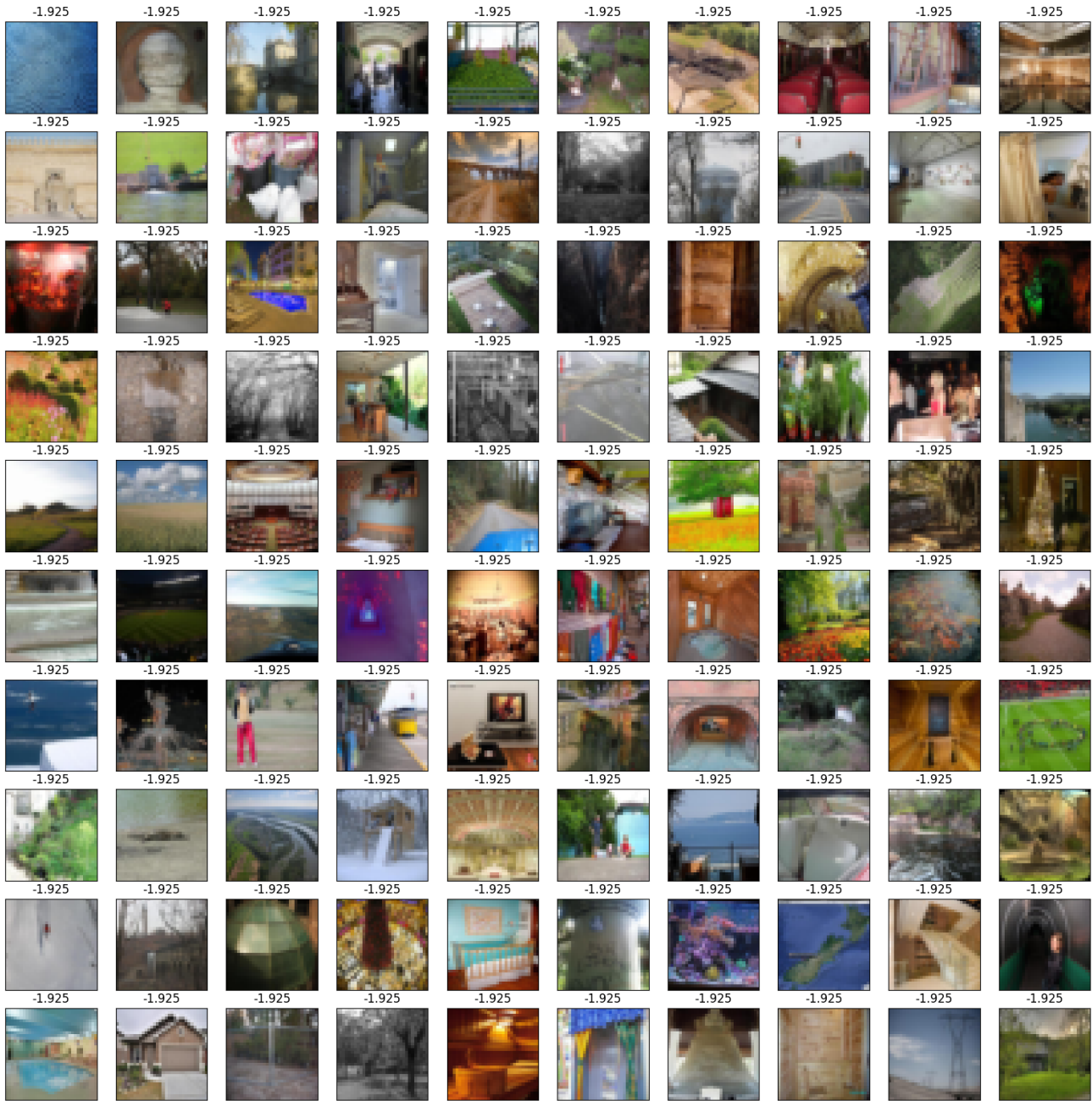


Figure 10: The set of top-100 images from the Places 365 data set which Hopfield Boosting recognized as in-distribution. The image captions show $s(\boldsymbol{\xi})$ of the respective image below the caption.

Figure 11: The set of top-100 images from the Places 365 data set which Hopfield Boosting recognized as out-of-distribution. The image captions show $s(\boldsymbol{\xi})$ of the respective image below the caption.

### F.5. Results on Noticeably Different Data Sets

The choice of additional data sets should not be driven by a desire to showcase good performance; rather, we suggest opting for data that highlights weaknesses, as it holds the potential to drive investigations and uncover novel insights. Simple toy data is preferable due to its typically clearer and more intuitive characteristics compared to complex natural image data. In alignment with these considerations, the following data sets captivated our interest: iCartoonFace (Zheng et al., 2020), Four Shapes (smeschke, 2018), and Retail Product Checkout (RPC) (Wei et al., 2022b). In Figure 12, we show random samples from these data sets to demonstrate the noticeable differences compared to CIFAR-10.

Table 7: Comparison between EBO-OE (Liu et al., 2020) and our version. ↓ indicates "lower is better" and ↑ indicates "higher is better".

| | Hopfield Boosting | | EBO-OE | |
|---|---|---|---|---|
| OOD Dataset | FPR95 ↓ | AUROC ↑ | FPR95 ↓ | AUROC ↑ |
| iCartoonFace | **0.60** | **99.57** | 4.01 | 98.94 |
| Four Shapes | **40.81** | **90.53** | 62.55 | 75.34 |
| RPC | **4.07** | **98.65** | 18.51 | 96.10 |



Figure 12: Random samples from three data sets, each noticeably different from CIFAR-10. First row: iCartoonFace; Second row: Four shapes; Third row: RPC.

In Table 7, we present some preliminary results using models trained with the respective method on CIFAR-10 as ID data set (as in Table 1). Results for comparison are presented for EBO only, as time constraints prevented experimenting with additional baseline methods. Although one would expect near-perfect results due to the evident disparities with CIFAR-10, Four Shapes (smeschke, 2018) and RPC (Wei et al., 2022b) seem to defy that expectation. Their results indicate a weakness in the capability to identify outliers robustly since many samples are classified as inliers. Only iCartoonFace (Zheng et al., 2020) is correctly detected as OOD, at least to a large degree. Interestingly, the weakness uncovered by this data is present in both methods, although more pronounced in EBO-OE. Therefore, we suspect that this specific behavior may be a general weakness when training OOD detectors using OE, an aspect we plan to investigate further in our future work.

### F.6. Runtime Considerations for Inference

When using Hopfield Boosting in inference, an additional inference step is needed to check whether a given sample is ID or OOD. Namely, to obtain the score (Equation (13)) of a query sample $\xi^{\mathcal{D}}$, Hopfield Boosting computes the dot product similarity of the embedding obtained from $\xi = \phi(\xi^{\mathcal{D}})$ to all samples in the Hopfield memories $X$ and $O$. In our experiments, $X$ contains the full in-distribution data set (50,000 samples) and $O$ contains a subset of the AUX data set of equal size. We investigate the computational overhead of computing the dot-product similarity to 100,000 samples in relation to the computational load of the encoder. For this, we feed 100 batches of size 1024 to an encoder (1) without using the score and (2) with using the score, measure the runtimes per batch, and compute the mean and standard deviation. We conduct this experiment with four different encoders on an NVIDIA Titan V GPU. The results are shown in Figure 13 and Table 8. One can see that, especially for larger models, the computational overhead of determining the score is very moderate in comparison.
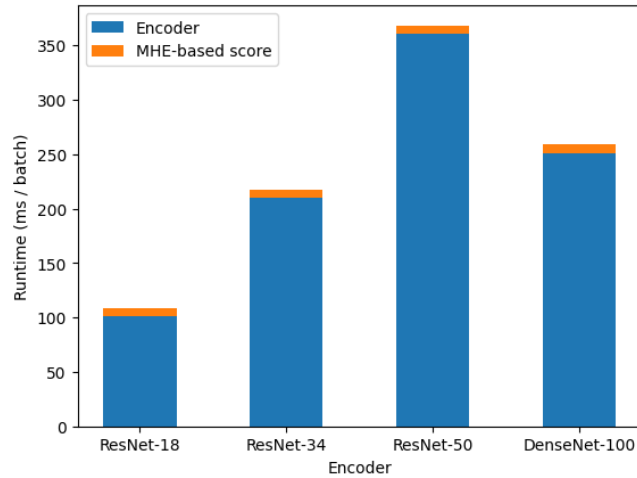
Figure 13: Mean inference runtimes for Hopfield Boosting on four different encoders on an NVIDIA Titan V GPU. We plot the contributions to the total runtime of the encoder and the MHE-based score (Equation (13)) separately. The evaluation shows that the score computation adds a negligible amount of computational overhead to the total runtime.

Table 8: Inference runtimes for Hopfield Boosting with four different encoders on an NVIDIA Titan V GPU. We compare the runtime of the encoder only and the runtime of the encoder with the MHE-based score computation (Equation (13)) combined.

| Encoder | Time encoder (ms / batch) | Time encoder + score (ms / batch) | Rel. overhead (%) |
|---|---|---|---|
| ResNet-18 | $100.93^{\pm0.24}$ | $108.50^{\pm0.19}$ | 7.50 |
| ResNet-34 | $209.80^{\pm0.40}$ | $217.33^{\pm0.51}$ | 3.59 |
| ResNet-50 | $360.93^{\pm1.51}$ | $368.17^{\pm0.62}$ | 2.01 |
| Densenet-100 | $251.24^{\pm1.36}$ | $258.82^{\pm0.84}$ | 3.02 |

Table 9: OOD detection performance on CIFAR-10. We compare results from Hopfield Boosting, PALM (Lu et al., 2024), NPOS (Tao et al., 2023), SSD+ (Sehwag et al., 2021), ASH (Djurisic et al., 2023), GEN (Liu et al., 2023), EBO (Liu et al., 2020), MaxLogit (Hendrycks et al., 2019a), and MSP (Hendrycks & Gimpel, 2017) on ResNet-18. ↓ indicates "lower is better" and ↑ "higher is better". All values in %. Standard deviations are estimated across five training runs.

| | | HB (ours) | PALM | NPOS | SSD+ | ASH | GEN | EBO | MaxLogit | MSP |
|---|---|---|---|---|---|---|---|---|---|---|
| SVHN | FPR95 ↓ | $\mathbf{0.23^{\pm 0.08}}$ | $1.24^{\pm 0.49}$ | $9.04^{\pm 1.13}$ | $3.05^{\pm 0.22}$ | $25.17^{\pm 9.55}$ | $33.26^{\pm 5.99}$ | $32.10^{\pm 6.41}$ | $33.27^{\pm 6.18}$ | $49.41^{\pm 3.77}$ |
| | AUROC ↑ | $99.57^{\pm 0.06}$ | $\mathbf{99.70^{\pm 0.12}}$ | $98.37^{\pm 0.23}$ | $99.41^{\pm 0.06}$ | $94.86^{\pm 2.09}$ | $93.53^{\pm 1.42}$ | $93.43^{\pm 1.60}$ | $93.29^{\pm 1.57}$ | $92.48^{\pm 0.93}$ |
| LSUN-Crop | FPR95 ↓ | $\mathbf{0.82^{\pm 0.17}}$ | $1.21^{\pm 0.27}$ | $5.52^{\pm 0.50}$ | $2.83^{\pm 1.10}$ | $13.13^{\pm 1.81}$ | $19.40^{\pm 2.22}$ | $17.25^{\pm 2.30}$ | $18.50^{\pm 2.24}$ | $38.32^{\pm 2.61}$ |
| | AUROC ↑ | $99.40^{\pm 0.04}$ | $\mathbf{99.65^{\pm 0.05}}$ | $98.97^{\pm 0.04}$ | $99.37^{\pm 0.16}$ | $97.33^{\pm 0.36}$ | $96.48^{\pm 0.46}$ | $96.73^{\pm 0.46}$ | $96.52^{\pm 0.47}$ | $94.37^{\pm 0.53}$ |
| LSUN-Resize | FPR95 ↓ | $\mathbf{0.00^{\pm 0.00}}$ | $27.01^{\pm 5.82}$ | $26.85^{\pm 3.14}$ | $34.30^{\pm 2.17}$ | $38.18^{\pm 5.78}$ | $31.50^{\pm 3.92}$ | $30.69^{\pm 4.03}$ | $31.64^{\pm 4.01}$ | $45.82^{\pm 3.48}$ |
| | AUROC ↑ | $\mathbf{99.98^{\pm 0.02}}$ | $95.41^{\pm 0.74}$ | $95.68^{\pm 0.36}$ | $94.78^{\pm 0.25}$ | $90.39^{\pm 2.00}$ | $94.04^{\pm 0.84}$ | $94.02^{\pm 0.86}$ | $93.90^{\pm 0.86}$ | $92.84^{\pm 0.80}$ |
| Textures | FPR95 ↓ | $\mathbf{0.16^{\pm 0.02}}$ | $17.32^{\pm 2.50}$ | $27.72^{\pm 2.55}$ | $21.20^{\pm 2.20}$ | $46.08^{\pm 6.22}$ | $44.62^{\pm 4.14}$ | $44.67^{\pm 4.46}$ | $44.97^{\pm 4.44}$ | $55.04^{\pm 2.86}$ |
| | AUROC ↑ | $\mathbf{99.84^{\pm 0.01}}$ | $96.82^{\pm 0.71}$ | $95.36^{\pm 0.35}$ | $96.46^{\pm 0.35}$ | $88.32^{\pm 2.08}$ | $90.12^{\pm 1.32}$ | $89.61^{\pm 1.50}$ | $89.56^{\pm 1.48}$ | $90.10^{\pm 0.92}$ |
| iSUN | FPR95 ↓ | $\mathbf{0.00^{\pm 0.00}}$ | $25.71^{\pm 4.83}$ | $26.90^{\pm 3.52}$ | $35.71^{\pm 2.27}$ | $42.41^{\pm 6.28}$ | $35.85^{\pm 4.05}$ | $34.99^{\pm 4.33}$ | $36.02^{\pm 4.18}$ | $49.10^{\pm 3.06}$ |
| | AUROC ↑ | $\mathbf{99.97^{\pm 0.02}}$ | $95.60^{\pm 0.65}$ | $95.74^{\pm 0.25}$ | $94.49^{\pm 0.25}$ | $89.06^{\pm 2.26}$ | $93.05^{\pm 0.84}$ | $92.99^{\pm 0.90}$ | $92.88^{\pm 0.90}$ | $91.99^{\pm 0.74}$ |
| Places 365 | FPR95 ↓ | $\mathbf{4.28^{\pm 0.23}}$ | $22.97^{\pm 2.17}$ | $32.62^{\pm 0.13}$ | $24.99^{\pm 1.21}$ | $48.03^{\pm 2.04}$ | $45.82^{\pm 1.07}$ | $44.87^{\pm 1.11}$ | $45.63^{\pm 1.26}$ | $57.58^{\pm 0.97}$ |
| | AUROC ↑ | $\mathbf{98.51^{\pm 0.10}}$ | $94.95^{\pm 0.53}$ | $93.76^{\pm 0.12}$ | $94.93^{\pm 0.22}$ | $85.65^{\pm 0.77}$ | $88.68^{\pm 0.28}$ | $88.53^{\pm 0.30}$ | $88.42^{\pm 0.29}$ | $88.06^{\pm 0.25}$ |
| Mean | FPR95 ↓ | $\mathbf{0.92}$ | 15.91 | 21.44 | 20.35 | 35.50 | 35.07 | 34.09 | 35.00 | 49.21 |
| | AUROC ↑ | $\mathbf{99.55}$ | 97.02 | 96.31 | 96.57 | 90.94 | 92.65 | 92.55 | 92.43 | 91.64 |
| Method type | | OE | Training | Training | Training | Post-hoc | Post-hoc | Post-hoc | Post-hoc | Post-hoc |
| Augmentations | | Weak | Strong | Strong | Strong | Weak | Weak | Weak | Weak | Weak |
| Auxiliary outlier data | | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

## F.7. Non-OE Baselines

To confirm the prevailing notion that OE methods can improve the OOD detection capability in general, we compare Hopfield Boosting to 3 training methods (Sehwag et al., 2021; Tao et al., 2023; Lu et al., 2024) and 5 post-hoc methods (Hendrycks & Gimpel, 2017; Hendrycks et al., 2019b; Liu et al., 2020; 2023; Djurisic et al., 2023). For all methods, we train a ResNet-18 on CIFAR-10. For Hopfield Boosting, we use the same training setup as described in section 5.1. For the post-hoc methods, we do not use the auxiliary outlier data. For the training methods, we use the training procedures described in the respective publications for 100 epochs. Notably, all training methods employ stronger augmentations than the OE or the post-hoc methods. The OE and post-hoc methods use the following augmentations (denoted as "Weak"):

1. RandomCrop (32x32), padding 4

2. RandomHorizontalFlip

The training methods use the following augmentations (denoted as "Strong"):

1. RandomResizedCrop (32x32), scale 0.2-1

2. RandomHorizontalFlip

3. ColorJitter applied with probability 0.8

4. RandomGrayscale applied with probability 0.2

Table 9 shows the results of the comparison of Hopfield Boosting to the post-hoc and training methods. Hopfield Boosting is better at OOD detection than all non-OE baselines on CIFAR-10 in terms of both mean AUROC and mean FPR95 by a large margin. Further, Hopfield Boosting achieves the best OOD detection on all OOD data sets in terms of FPR95 and AUROC, except for SVHN and LSUN-Crop, where PALM (Lu et al., 2024) shows better AUROC results. An interesting avenue for future work is to combine one of the non-OE based training methods with the OE method Hopfield Boosting.

## G. Informativeness of Sampling with High Boundary Scores

This section adopts and expands the arguments of Ming et al. (2022) on sampling with high boundary scores.

We assume the extracted features of a trained deep neural network to approximately equal a Gaussian mixture model with equal class priors:

$$p(\boldsymbol{\xi}) \;=\; \frac{1}{2}\mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) + \frac{1}{2}\mathcal{N}(\boldsymbol{\xi}; -\boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) \tag{125}$$

$$p_{\text{ID}}(\boldsymbol{\xi}) = p(\boldsymbol{\xi}|\text{ID}) \;=\; \mathcal{N}(\boldsymbol{\xi}; \boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) \tag{126}$$

$$p_{\text{AUX}}(\boldsymbol{\xi}) = p(\boldsymbol{\xi}|\text{AUX}) \;=\; \mathcal{N}(\boldsymbol{\xi}; -\boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) \tag{127}$$

Using the MHE and sufficient data from those distributions, we can estimate the densities $p(\boldsymbol{\xi})$, $p(\boldsymbol{\xi}|\text{ID})$ and $p(\boldsymbol{\xi}|\text{AUX})$.

**Lemma G.1.** *(see Lemma E.1 in Ming et al. (2022)) Assume the $M$ sampled data points $\boldsymbol{o}_i \sim p_{AUX}$ satisfy the following constraint on high boundary scores $\mathrm{E}_b(\boldsymbol{\xi})$*

$$\frac{-\sum_{i=1}^{M} \mathrm{E}_b(\boldsymbol{o}_i)}{M} \leq \epsilon \tag{128}$$

*Then they have*

$$\sum_{i=1}^{M} |2\boldsymbol{\mu}^T \boldsymbol{o}_i| \leq M\epsilon\sigma^2 \tag{129}$$

*Proof.* They first obtain the expression for $\mathrm{E}_b(\boldsymbol{\xi})$ under the Gaussian mixture model described above and can express $p(\text{AUX}|\boldsymbol{\xi})$ as

$$p(\text{AUX}|\boldsymbol{\xi}) \;=\; \frac{p(\boldsymbol{\xi}|\text{AUX})p(\text{AUX})}{p(\boldsymbol{\xi})} \tag{130}$$

$$=\; \frac{\frac{1}{2}p(\boldsymbol{\xi}|\text{AUX})}{\frac{1}{2}p(\boldsymbol{\xi}|\text{ID}) \;+\; \frac{1}{2}p(\boldsymbol{\xi}|\text{AUX})} \tag{131}$$

$$=\; \frac{(2\pi\sigma^2)^{-d/2}\exp(-\frac{1}{2\sigma^2}\|\boldsymbol{\xi}-\boldsymbol{\mu}\|_2^2)}{(2\pi\sigma^2)^{-d/2}\exp(-\frac{1}{2\sigma^2}\|\boldsymbol{\xi}+\boldsymbol{\mu}\|_2^2) \;+\; (2\pi\beta^{-1})^{-d/2}\exp(-\frac{1}{2\sigma^2}\|\boldsymbol{\xi}-\boldsymbol{\mu}\|_2^2)} \tag{132}$$

$$=\; \frac{1}{1 \;+\; \exp(-\frac{1}{2\sigma^2}(\|\boldsymbol{\xi}-\boldsymbol{\mu}\|_2^2 - \|\boldsymbol{\xi}+\boldsymbol{\mu}\|_2^2))} \tag{133}$$

When defining $f_{\text{AUX}}(\boldsymbol{\xi}) \;=\; \frac{1}{2\sigma^2}(\|\boldsymbol{\xi}-\boldsymbol{\mu}\|_2^2 - \|\boldsymbol{\xi}+\boldsymbol{\mu}\|_2^2)$ such that $p(\text{AUX}|\boldsymbol{\xi}) \;=\; \sigma(f_{\text{AUX}}(\boldsymbol{\xi})) \;=\; \frac{1}{1 \;+\; \exp(-f_{\text{AUX}}(\boldsymbol{\xi}))}$, they define $\mathrm{E}_b$ as follows:

$$\mathrm{E}_b(\boldsymbol{\xi}) \;=\; -|f_{\text{AUX}}(\boldsymbol{\xi})| \tag{134}$$

$$=\; -\frac{1}{2\sigma^2}|\,\|\boldsymbol{\xi}-\boldsymbol{\mu}\|_2^2 - \|\boldsymbol{\xi}+\boldsymbol{\mu}\|_2^2\,| \tag{135}$$

$$=\; -\frac{1}{2\sigma^2}|\,\boldsymbol{\xi}^T\boldsymbol{\xi} - 2\boldsymbol{\mu}^T\boldsymbol{\xi} + \boldsymbol{\mu}^T\boldsymbol{\mu} - (\boldsymbol{\xi}^T\boldsymbol{\xi} + 2\boldsymbol{\mu}^T\boldsymbol{\xi} + \boldsymbol{\mu}^T\boldsymbol{\mu})\,| \tag{136}$$

$$=\; -\frac{|2\boldsymbol{\mu}^T\boldsymbol{\xi}|}{\sigma^2} \tag{137}$$

Therefore, the constraint in Equation (129) is translated to

$$\sum_{i=1}^{M} |2\boldsymbol{\mu}^T \boldsymbol{o}_i| \leq M\epsilon\sigma^2 \tag{138}$$

$\square$

As $\max_{i \in M} |\boldsymbol{\mu}^T \boldsymbol{o}_i| \leq \sum_{i=1}^{M} |\boldsymbol{\mu}^T \boldsymbol{o}_i|$ given a fixed $M$, the selected samples can be seen as generated from $p_{\text{AUX}}$ with the constraint that all samples lie within the two hyperplanes in Equation (138).

**Parameter estimation.** Now they show the benefit of such constraint in controlling the sample complexity. Assume the signal/noise ratio is large: $\frac{||\boldsymbol{\mu}||}{\sigma} = r \gg 1$, and $\epsilon \leq 1$ is some constant.

Assume the classifier is given by

$$\boldsymbol{\theta} = \frac{1}{N+M}\left(\sum_{i=1}^{M} \boldsymbol{x}_i - \sum_{i=1}^{N} \boldsymbol{o}_i\right) \tag{139}$$

where $\boldsymbol{o}_i \sim p_{\text{AUX}}$ and $\boldsymbol{x}_i \sim p_{\text{ID}}$. One can decompose $\boldsymbol{\theta}$. Assuming $M = N$:

$$\boldsymbol{\theta} = \boldsymbol{\mu} + \frac{1}{2}\boldsymbol{\eta} + \frac{1}{2}\boldsymbol{\omega} \tag{140}$$

$$\boldsymbol{\eta} = \frac{1}{N}\left(\sum_{i=1}^{N} \boldsymbol{x}_i\right) - \boldsymbol{\mu} \tag{141}$$

$$\boldsymbol{\omega} = \frac{1}{N}\left(\sum_{i=1}^{M} - \boldsymbol{o}_i\right) - \boldsymbol{\mu} \tag{142}$$

We would now like to determine the distributions of the random variables $||\boldsymbol{\eta}||_2^2$ and $\boldsymbol{\mu}^T \boldsymbol{\eta}$

$$||\boldsymbol{\eta}||_2^2 = \sum_{i=1}^{d} \eta_i^2 \tag{143}$$

$$\eta_i \sim \mathcal{N}(0, \frac{\sigma^2}{N}) \tag{144}$$

$$\frac{\sqrt{N}}{\sigma}\eta_i \sim \mathcal{N}(0, 1) \tag{145}$$

$$(\frac{\sqrt{N}}{\sigma}\eta_i)^2 \sim \chi_1^2 \tag{146}$$

Therefore, for $||\boldsymbol{\eta}||_2^2$ we have

$$\frac{N}{\sigma^2}||\boldsymbol{\eta}||_2^2 = \sum_{i=1}^{d}(\frac{\sqrt{N}}{\sigma}\eta_i)^2 \sim \chi_d^2 \tag{147}$$

Now we would like to determine the distribution of $\boldsymbol{\mu}^T \boldsymbol{\eta}$:

$$\boldsymbol{\mu}^T \boldsymbol{\eta} = \sum_{i=1}^{d} \mu_i \, \eta_i \tag{148}$$

$$\mu_i \, \eta_i \sim \mathcal{N}(0, \frac{\sigma^2 \mu_i^2}{N}) \tag{149}$$

$$\sum_{i=1}^{d} \mu_i \, \eta_i \sim \mathcal{N}(0, \sum_{i=1}^{d} \frac{\sigma^2 \mu_i^2}{N}) \tag{150}$$

$$\sum_{i=1}^{d} \mu_i \, \eta_i \sim \mathcal{N}(0, \frac{\sigma^2}{N} \sum_{i=1}^{d} \mu_i^2) \tag{151}$$

$$\frac{\boldsymbol{\mu}^T \boldsymbol{\eta}}{||\boldsymbol{\mu}||} \sim \mathcal{N}(0, \frac{\sigma^2}{N}) \tag{152}$$

**Concentration bounds.** They now develop concentration bounds for $||\boldsymbol{\eta}||_2^2$ and $\boldsymbol{\mu}^T \boldsymbol{\eta}$. First, we look at $||\boldsymbol{\eta}||_2^2$. A concentration bound for $\chi_d^2$ is:

$$\mathbb{P}(X - d \geq 2\sqrt{dx} + 2x) \leq \exp(-x) \tag{153}$$

By assuming $x = \frac{d}{8\sigma^2}$ we obtain

$$\mathbb{P}(X - d \geq 2\sqrt{d\frac{d}{8\sigma^2}} + 2\frac{d}{8\sigma^2}) \leq \exp(-\frac{d}{8\sigma^2}) \tag{154}$$

$$\mathbb{P}(X \geq \frac{d}{\sqrt{2}\sigma} + \frac{d}{4\sigma^2} + d) \leq \exp(-\frac{d}{8\sigma^2}) \tag{155}$$

$$\mathbb{P}(\frac{N}{\sigma^2}||\boldsymbol{\eta}||_2^2 \geq \frac{d}{\sqrt{2}\sigma} + \frac{d}{4\sigma^2} + d) \leq \exp(-\frac{d}{8\sigma^2}) \tag{156}$$

$$\mathbb{P}(||\boldsymbol{\eta}||_2^2 \geq \frac{\sigma^2}{N}(\frac{d}{\sqrt{2}\sigma} + \frac{d}{4\sigma^2} + d)) \leq \exp(-\frac{d}{8\sigma^2}) \tag{157}$$

If $d \geq 2$ we have that[3]

$$\frac{d}{\sqrt{2}\sigma} + \frac{d}{4\sigma^2} > \frac{1}{\sigma} \tag{158}$$

and thus, the above bound can be simplified when assuming $d \geq 2$ as follows:

$$\mathbb{P}(||\boldsymbol{\eta}||_2^2 \geq \frac{\sigma^2}{N}(\frac{1}{\sigma} + d)) \leq \exp(-\frac{d}{8\sigma^2}) \tag{159}$$

For $||\boldsymbol{\omega}||_2^2$, since all $\boldsymbol{o}_i$ is drawn i.i.d. from $p_{\text{AUX}}$, under the constraint in Equation (138), the distribution of $\boldsymbol{\omega}$ can be seen as a truncated distribution of $\boldsymbol{\eta}$. Thus, with some finite positive constant $c$, we have

$$\mathbb{P}(||\boldsymbol{\omega}||_2^2 \geq \frac{\sigma^2}{N}(d + \frac{1}{\sigma})) \leq c\mathbb{P}(||\boldsymbol{\eta}||_2^2 \geq \frac{\sigma^2}{N}(d + \frac{1}{\sigma})) \leq c\exp(-\frac{d}{8\sigma^2}) \tag{160}$$

---

[3]Strictly, the bound is valid for $d > \sqrt{2}$

Now, we develop a bound for $\boldsymbol{\mu}^T\boldsymbol{\eta}$. A concentration bound for $\mathcal{N}(\mu, \sigma^2)$ is

$$\mathbb{P}(X - \mu \geq t) \leq \exp(\frac{-t^2}{2\sigma^2}) \qquad (161)$$

By applying $\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \sim \mathcal{N}(0, \frac{\sigma^2}{N})$ to the above bound we obtain

$$\mathbb{P}(\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \geq t) \leq \exp(\frac{-t^2 N}{2\sigma^2}) \qquad (162)$$

Assuming $t = (\sigma||\boldsymbol{\mu}||)^{1/2}$ we obtain

$$\mathbb{P}(\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \geq (\sigma||\boldsymbol{\mu}||)^{1/2}) \leq \exp(\frac{-(\sigma||\boldsymbol{\mu}||)N}{2\sigma^2}) \qquad (163)$$

$$\mathbb{P}(\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \geq (\sigma||\boldsymbol{\mu}||)^{1/2}) \leq \exp(\frac{-||\boldsymbol{\mu}||N}{2\sigma}) \qquad (164)$$

Due to symmetry, we have

$$\mathbb{P}(-\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \leq -(\sigma||\boldsymbol{\mu}||)^{1/2}) \leq \exp(\frac{-||\boldsymbol{\mu}||N}{2\sigma}) \qquad (165)$$

$$\mathbb{P}(-\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \leq -(\sigma||\boldsymbol{\mu}||)^{1/2}) + \mathbb{P}(\frac{\boldsymbol{\mu}^T\boldsymbol{\eta}}{||\boldsymbol{\mu}||} \geq (\sigma||\boldsymbol{\mu}||)^{1/2}) \leq 2\exp(\frac{-||\boldsymbol{\mu}||N}{2\sigma}) \qquad (166)$$

We can rewrite the above bound using the absolute value function.

$$\mathbb{P}(\frac{|\boldsymbol{\mu}^T\boldsymbol{\eta}|}{||\boldsymbol{\mu}||} \geq (\sigma||\boldsymbol{\mu}||)^{1/2}) \leq 2\exp(\frac{-||\boldsymbol{\mu}||N}{2\sigma}) \qquad (167)$$

**Benefit of high boundary scores.** We will now show why sampling with high boundary scores is beneficial. Recall the results from Equations (138) and (142):

$$\sum_{i=1}^{M} |2\boldsymbol{\mu}^T\boldsymbol{o}_i| \leq M\epsilon\sigma^2 \qquad (168)$$

$$\boldsymbol{\omega} = \frac{1}{M}(-\sum_{i=1}^{M} \boldsymbol{o}_i) - \boldsymbol{\mu} \qquad (169)$$

The triangle inequality is

$$|a + b| \leq |a| + |b| \qquad (170)$$

$$|a + (-b)| \leq |a| + |b| \qquad (171)$$

Using the two facts above and the triangle inequality we can bound $|\boldsymbol{\mu}^T\boldsymbol{\omega}|$:

$$\frac{1}{M}|\sum_{i=1}^{M}\boldsymbol{\mu}^T\boldsymbol{o}_i| \leq \frac{\sigma^2\epsilon}{2} \tag{172}$$

$$\frac{1}{M}|-\sum_{i=1}^{M}\boldsymbol{\mu}^T\boldsymbol{o}_i| \leq \frac{\sigma^2\epsilon}{2} \tag{173}$$

$$\frac{1}{M}|-\sum_{i=1}^{M}\boldsymbol{\mu}^T\boldsymbol{o}_i| + ||\boldsymbol{\mu}||_2^2 \leq \frac{\sigma^2\epsilon}{2} + ||\boldsymbol{\mu}||_2^2 \tag{174}$$

$$\frac{1}{M}|-\sum_{i=1}^{M}\boldsymbol{\mu}^T\boldsymbol{o}_i - \boldsymbol{\mu}^T\boldsymbol{\mu}| \leq \frac{\sigma^2\epsilon}{2} + ||\boldsymbol{\mu}||_2^2 \tag{175}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\omega}| \leq ||\boldsymbol{\mu}||_2^2 + \frac{\sigma^2\epsilon}{2} \tag{176}$$

**Developing a lower bound.** Let

$$||\boldsymbol{\eta}||_2^2 \leq \frac{\sigma^2}{N}(d + \frac{1}{\sigma}) \tag{177}$$

$$||\boldsymbol{\omega}||_2^2 \leq \frac{\sigma^2}{N}(d + \frac{1}{\sigma}) \tag{178}$$

$$\frac{|\boldsymbol{\mu}^T\boldsymbol{\eta}|}{||\boldsymbol{\mu}||} \leq (\sigma||\boldsymbol{\mu}||)^{1/2} \tag{179}$$

hold simultaneously. The probability of this happening can be bounded as follows: We define $T$ and its complement $\bar{T}$:

$$T = \{||\boldsymbol{\eta}||_2^2 \leq \frac{\sigma^2}{N}(d + \frac{1}{\sigma})\} \cap \{||\boldsymbol{\omega}||_2^2 \leq \frac{\sigma^2}{N}(d + \frac{1}{\sigma})\} \cap \{\frac{|\boldsymbol{\mu}^T\boldsymbol{\eta}|}{||\boldsymbol{\mu}||} \leq (\sigma||\boldsymbol{\mu}||)^{1/2}\} \tag{180}$$

$$\bar{T} = \{||\boldsymbol{\eta}||_2^2 > \frac{\sigma^2}{N}(d + \frac{1}{\sigma})\} \cup \{||\boldsymbol{\omega}||_2^2 > \frac{\sigma^2}{N}(d + \frac{1}{\sigma})\} \cup \{\frac{|\boldsymbol{\mu}^T\boldsymbol{\eta}|}{||\boldsymbol{\mu}||} > (\sigma||\boldsymbol{\mu}||)^{1/2}\} \tag{181}$$

With $\mathbb{P}(T) + \mathbb{P}(\bar{T}) = 1$. The probability $\mathbb{P}(\bar{T})$ can be bounded using Boole's inequality and the results in Equations (159), (160) and (167):

$$\mathbb{P}(\bar{T}) \leq \exp(-d/8\sigma^2) + c\exp(-d/8\sigma^2) + 2\exp(\frac{-||\mu||N}{2\sigma}) \tag{182}$$

$$\mathbb{P}(\bar{T}) \leq (1 + c)\exp(-d/8\sigma^2) + 2\exp(\frac{-||\mu||N}{2\sigma}) \tag{183}$$

Further, we can bound the probability $\mathbb{P}(T)$:

$$\mathbb{P}(\bar{T}) \leq (1 + c)\exp(-d/8\sigma^2) + 2\exp(\frac{-||\mu||N}{2\sigma}) \tag{184}$$

$$1 - \mathbb{P}(T) \leq (1 + c)\exp(-d/8\sigma^2) + 2\exp(\frac{-||\mu||N}{2\sigma}) \tag{185}$$

$$\mathbb{P}(T) \geq 1 - (1 + c)\exp(-d/8\sigma^2) - 2\exp(\frac{-||\mu||N}{2\sigma}) \tag{186}$$

Therefore, the probability of the assumptions in Equations (177), (178), and (179) occuring simultneously is at least $1 - (1 + c) \exp(-d/8\sigma^2) - 2 \exp(\frac{-||\boldsymbol{\mu}||N}{2\sigma})$.

By using the triangle inequality, Equation (140) and the Assumptions (177) and (178) we can bound $||\boldsymbol{\theta}||_2^2$:

$$||\boldsymbol{\theta}||_2^2 = ||\ \boldsymbol{\mu}\ +\ \frac{1}{2}\ \boldsymbol{\eta}\ +\ \frac{1}{2}\ \boldsymbol{\omega}||_2^2 \tag{187}$$

$$||\boldsymbol{\theta}||_2^2 \leq ||\boldsymbol{\mu}||_2^2 + ||\frac{1}{2}\ \boldsymbol{\eta}||_2^2 + ||\frac{1}{2}\ \boldsymbol{\omega}||_2^2 \tag{188}$$

$$||\boldsymbol{\theta}||_2^2 \leq ||\boldsymbol{\mu}||_2^2 + \frac{1}{4}||\boldsymbol{\eta}||_2^2 + \frac{1}{4}||\boldsymbol{\omega}||_2^2 \tag{189}$$

$$||\boldsymbol{\theta}||_2^2 \leq ||\boldsymbol{\mu}||_2^2 + \frac{1}{2}\frac{\sigma^2}{N}(d + \frac{1}{\sigma}) \tag{190}$$

$$||\boldsymbol{\theta}||_2^2 \leq ||\boldsymbol{\mu}||_2^2 + \frac{\sigma^2}{N}(d + \frac{1}{\sigma}) \tag{191}$$

The reverse triangle inequality is defined as

$$|x - y| \geq \big||x| - |y|\big| \tag{192}$$

$$|x - (-y)| \geq \big||x| - |y|\big| \tag{193}$$

Using the reverse triangle inequality, Equations (140), (176) and Assumption (179) we have that

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| = |\boldsymbol{\mu}^T\boldsymbol{\mu}\ +\ \frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\eta}\ +\ \frac{1}{2}\ \boldsymbol{\mu}^T\boldsymbol{\omega}| \tag{194}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big||\boldsymbol{\mu}^T\boldsymbol{\mu}|\ -\ |\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\eta}|\ -\ |\frac{1}{2}\ \boldsymbol{\mu}^T\boldsymbol{\omega}|\big| \tag{195}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big|||\boldsymbol{\mu}||_2^2\ -\ \frac{1}{2}\sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{1}{2}||\boldsymbol{\mu}||_2^2\ -\ \frac{1}{2}\frac{\sigma^2\epsilon}{2}\big| \tag{196}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big|\frac{1}{2}||\boldsymbol{\mu}||_2^2\ -\ \frac{1}{2}\sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{1}{2}\frac{\sigma^2\epsilon}{2}\big| \tag{197}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big|\frac{1}{2}(||\boldsymbol{\mu}||_2^2\ -\ \sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{\sigma^2\epsilon}{2})\big| \tag{198}$$

We have assumed that the signal/noise ratio is large: $\frac{||\boldsymbol{\mu}||}{\sigma} = r \gg 1$. Thus, we can drop the absolute value, because we assume that the term inside the $||$ is larger than zero:

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big|\frac{1}{2}(||\boldsymbol{\mu}||_2^2\ -\ \frac{1}{r}||\boldsymbol{\mu}||^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{||\boldsymbol{\mu}||_2^2\epsilon}{2r^2})\big| \tag{199}$$

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \big|(1 - \frac{1}{r} - \frac{\epsilon}{2r^2})\frac{1}{2}(||\boldsymbol{\mu}||_2^2)\big| \tag{200}$$

We have

$$(1 - \frac{1}{r} - \frac{\epsilon}{2r^2}) \geq 0 \tag{201}$$

if $r \geq 1.36602540378443\ldots$ and $\epsilon \leq 1$, and therefore

$$|\boldsymbol{\mu}^T\boldsymbol{\theta}| \geq \frac{1}{2}(||\boldsymbol{\mu}||_2^2 - \sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{\sigma^2\epsilon}{2}) \tag{202}$$

Because of Equation (191) and the fact that if $x \leq y$ and $\operatorname{sgn}(x) = \operatorname{sgn}(y)$ then $x^{-1} \geq y^{-1}$ we have

$$\frac{1}{||\boldsymbol{\theta}||} \geq \frac{1}{\sqrt{||\boldsymbol{\mu}||_2^2 + \frac{\sigma^2}{N}(d + \frac{1}{\sigma})}} \tag{203}$$

We can combine the Equations (202) and (203) to give a single bound:

$$\frac{|\boldsymbol{\mu}^T\boldsymbol{\theta}|}{||\boldsymbol{\theta}||} \geq \frac{||\boldsymbol{\mu}||_2^2 - \sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{\sigma^2\epsilon}{2}}{2\sqrt{||\boldsymbol{\mu}||_2^2 + \frac{\sigma^2}{N}(d + \frac{1}{\sigma})}} \tag{204}$$

we define $\boldsymbol{\theta}$ such that $\boldsymbol{\mu}^T\boldsymbol{\theta} > 0$ and thus

$$\frac{\boldsymbol{\mu}^T\boldsymbol{\theta}}{||\boldsymbol{\theta}||} \geq \frac{||\boldsymbol{\mu}||_2^2 - \sigma^{1/2}||\boldsymbol{\mu}||^{3/2} - \frac{\sigma^2\epsilon}{2}}{2\sqrt{||\boldsymbol{\mu}||_2^2 + \frac{\sigma^2}{N}(d + \frac{1}{\sigma})}} \tag{205}$$

The false negative rate $\mathrm{FNR}(\boldsymbol{\theta})$ and false positive rate $\mathrm{FPR}(\boldsymbol{\theta})$ are

$$\mathrm{FNR}(\boldsymbol{\theta}) = \int_{-\infty}^{0} \mathcal{N}(x; \frac{\boldsymbol{\mu}^T\boldsymbol{\theta}}{||\boldsymbol{\theta}||}, \sigma^2)\,\mathrm{d}x \tag{206}$$

$$\mathrm{FPR}(\boldsymbol{\theta}) = \int_{0}^{\infty} \mathcal{N}(x; \frac{-\boldsymbol{\mu}^T\boldsymbol{\theta}}{||\boldsymbol{\theta}||}, \sigma^2)\,\mathrm{d}x \tag{207}$$

As $\mathcal{N}(x; \mu, \sigma^2) = \mathcal{N}(-x; -\mu, \sigma^2)$, we have $\mathrm{FNR}(\boldsymbol{\theta}) = \mathrm{FPR}(\boldsymbol{\theta})$. From Equation (205) we can see that as $\epsilon$ decreases, the lower bound of $\frac{\boldsymbol{\mu}^T\boldsymbol{\theta}}{||\boldsymbol{\theta}||}$ will increase. Thus, the mean of the Gaussian distribution in Equation (206) will increase and therefore, the false negative rate will decrease, which shows the benefit of sampling with high boundary scores.