



HAL
open science

Provably Safeguarding a Classifier from OOD and Adversarial Samples: an Extreme Value Theory Approach

Nicolas Atienza, Christophe Labreuche, Johanne Cohen, Michèle Sebag

► **To cite this version:**

Nicolas Atienza, Christophe Labreuche, Johanne Cohen, Michèle Sebag. Provably Safeguarding a Classifier from OOD and Adversarial Samples: an Extreme Value Theory Approach. ICLR 2025 - The Thirteenth International Conference on Learning Representations, Apr 2025, Singapore (SG), Singapore. hal-04922382

HAL Id: hal-04922382

<https://inria.hal.science/hal-04922382v1>

Submitted on 30 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

PROVABLY SAFEGUARDING A CLASSIFIER FROM OOD AND ADVERSARIAL SAMPLES: AN EXTREME VALUE THEORY APPROACH

preprint

Nicolas Atienza
Thales cortAix-Labs
Industrial AI Laboratory SINCLAIR
Paris-Saclay University
Palaiseau, France

Christophe Labreuche
Thales cortAix-Labs
Industrial AI Laboratory SINCLAIR
Palaiseau, France

Johanne Cohen, Michèle Sebag
LISN CNRS-INRIA
Paris-Saclay University
Saclay, France

ABSTRACT

This paper introduces a novel method, *Sample-efficient Probabilistic Detection using Extreme Value Theory* (SPADE), which transforms a classifier into an abstaining classifier, offering provable protection against out-of-distribution and adversarial samples. The approach is based on a Generalized Extreme Value (GEV) model of the training distribution in the classifier’s latent space, enabling the formal characterization of OOD samples. Interestingly, under mild assumptions, the GEV model also allows for formally characterizing adversarial samples. The abstaining classifier, which rejects samples based on their assessment by the GEV model, provably avoids OOD and adversarial samples. The empirical validation of the approach, conducted on various neural architectures (ResNet, VGG, and Vision Transformer) and medium and large-sized datasets (CIFAR-10, CIFAR-100, and ImageNet), demonstrates its frugality, stability, and efficiency compared to the state of the art.

1 Introduction

A key challenge in deploying learned models in real-world settings is managing out-of-distribution (OOD) samples. When a learned model encounters data that deviates from the training distribution, it can lead to failures with significant consequences, particularly in high-stakes applications such as medical diagnosis, autonomous driving or risk analysis [1, 2]. The ultimate aim in machine learning is to achieve OOD generalization, where the model encapsulates the core concept with sufficient accuracy to effectively handle atypical but real samples [3]. A step towards OOD generalization is OOD detection, which equips the learned model with the ability to recognize atypical samples and refrain from making risky predictions. OOD detection is approached from several directions, including methods based on classification [4, 5, 6], reconstruction metrics [7, 8], density estimation [9, 10, 11] and distance-based estimation [12, 13, 14, 15] (more in Section 6).

OOD detection is complicated by the fact that, to the best of our knowledge, no universally accepted definition exists for what qualifies as an OOD sample. The boundaries between in-distribution and OOD data are inherently ambiguous and different domain experts may classify the same sample differently based on their understanding and experience [16]. Consequently, the validation of OOD detection methods relies heavily on experimental studies using well-curated datasets, such as near and far OOD datasets [2].

It is worth noting that human experts and models tend to make different decisions regarding both OOD and adversarial samples [17], albeit in distinct ways. Experts typically recognize an OOD sample as belonging to a given class, despite its atypicality, while the model assigns it to a random class. Conversely, experts perceive an adversarial example as typical of a specific class, yet the model confidently misclassifies it into a different class.

The approach presented in this paper, referred to as *Sample-efficient Probabilistic Detection using Extreme value theory* (SPADE) and inspired by distance-based approaches [13], introduces an original model of the training distribution relative to a learned model (hereafter the teacher). Specifically, the distances between samples in the teacher’s latent space are modelled using the Extreme Value Theory (EVT) [18]. This model provides a sound and robust test for detecting and rejecting OOD samples. Most interestingly, under mild assumptions this test also provably rejects adversarial examples with high probability, subject to a bound on their perturbation amplitude.

The contributions of the proposed approach are fourfold: i) it introduces a formal definition of OOD samples relative to a teacher and its latent representation; ii) this definition leads to a statistically frugal OOD test based on EVT first principles; iii) this test operationally rejects OOD and adversarial samples with provable guarantees; iv) the effectiveness of the approach is experimentally and successfully demonstrated against strong baselines for learned models with different architectures [19, 20, 15].

The paper is organized as follows. Section 2 outlines the formal background of distance-based OOD detection and introduces extreme value theory. Section 3 gives an overview of the SPADE approach and its formal analysis. Sections 4 and 5 respectively detail the experimental setup and the experiments conducted to validate SPADE against state-of-the-art methods. Section 6 discusses the contributions within the context of related work, and the paper concludes with perspectives for future research.

Notations. Let $\mathcal{D} = \{(\mathbf{x}_1, y_1) \dots (\mathbf{x}_n, y_n)\}$ denote the training set, iid drawn after the joint distribution $P_{\mathcal{X}, \mathcal{Y}}$, with \mathcal{X} the instance space and $\mathcal{Y} = \{1, \dots, n_c\}$ the set of classes. The trained teacher f is expressed as $f = g \circ h$, where h is the embedding in the latent space ($h : \mathcal{X} \mapsto \mathbb{R}^d$) and g is the mapping used to make decisions based on the latent representation.

2 Formal Background

This section describes the main concepts defined by [3] in the context of distance-based OOD generalization, and briefly introduces extreme value theory for completeness.

2.1 Properties of Latent In-Distribution

As mentioned above, the complexity of OOD characterization corresponds to the highly complex and diverse nature of real-world data [21]. In the literature [22, 23, 24, 25, 26] OOD characterization often involves subjective assessment ("data that appear noticeably different from in-distribution to human observers"). A common definition states that the OOD is different from the known in-distribution, e.g., the training distribution, although this definition does not capture the specifics of OOD: OOD is not merely any distribution that differs from ID.

On the other hand, as noted by [12, 13], the representation in latent space of the in-distribution (hereafter ID) presents distinct characteristics, such as being formed of compact and well-separated clusters. [3] formalize these properties in terms of variation and informativeness of the latent representation:

Definition 1 ([3]). *The variation of embedding h across a finite distribution \mathcal{D} , noted $\mathcal{V}_\rho(h, \mathcal{D})$, is defined as the maximum diameter over all classes c of the ball containing distribution $h(\mathbf{x})$ for $(\mathbf{x}, y) \in \mathcal{D}$ and $y = c$:*

$$\mathcal{V}(h, \mathcal{D}) = \max_{c \in \mathcal{Y}} \sup_{\substack{(\mathbf{x}, c) \in \mathcal{D} \\ (\mathbf{x}', c) \in \mathcal{D}}} \|h(\mathbf{x}) - h(\mathbf{x}')\| \quad (1)$$

where the distance $\|h(\mathbf{x}) - h(\mathbf{x}')\|$ is usually set to L_2 distance¹. Embedding h is said η -invariant across \mathcal{D} if $\mathcal{V}(h, \mathcal{D}) < \eta$.

The variation of embedding h thus measures the maximum thickness and width of the latent manifold containing the image of (samples in) a class.

Definition 2 ([3]). *The informativeness of embedding h across a finite distribution \mathcal{D} , noted $\mathcal{I}(h, \mathcal{D})$, is defined as the average over all pairs (c, c') of distinct classes, of the minimum distance between $h(\mathbf{x})$ and $h(\mathbf{x}')$ for (\mathbf{x}, y) and (\mathbf{x}', y') in \mathcal{D} , with $y = c \neq y' = c'$:*

$$\mathcal{I}(h, \mathcal{D}) = \frac{1}{n_c(n_c - 1)} \sum_{c \neq c' \in \mathcal{Y}} \min_{\substack{(\mathbf{x}, y) \in \mathcal{D}, (\mathbf{x}', y') \in \mathcal{D} \\ y = c \neq y' = c'}} \|h(\mathbf{x}) - h(\mathbf{x}')\| \quad (2)$$

¹The Kullback Leibler divergence is considered when embedding h is a probabilistic one. See [3] for more detail.

where the distance is usually set to L_2 distance (see footnote 1), and n_c denotes the number of classes. The latent embedding h is said δ -informative across \mathcal{D} if $\mathcal{I}(h, \mathcal{D}) > \delta$.

It is noted that informativeness and variation are closely related to the *compactness* and *dispersion* metrics introduced and optimized in CIDER to achieve OOD generalization [15].

These definitions are defined in terms of supremum over classes, raising the challenge of their statistical stability and algorithmic exploitation (e.g., through setting thresholds). The approach presented here addresses this challenge by exploiting extreme value theory, as described below and referring the reader to [27] for a comprehensive introduction.

2.2 Extreme Value Theory

Dating back to [18, 28], Extreme Value Theory (EVT) focuses on modeling and understanding the tail behavior of distributions. EVT is based on the premise that, under mild assumptions, the distributions of extreme events converge to a common form, even if their original distributions differ. For instance, while the distributions of seismic intensities and the heights of rogue waves – factors that respectively influence the design of buildings and oil rigs – may differ, their extreme values (maxima) are governed by the same class of distributions. This limiting distribution is known as the *Extreme Value Distribution* (EVD):

Definition 3 (Extreme Value Distribution (EVD) [18]). *Let Z be a random variable over the real-valued space \mathbb{R} . Let $Z^{(\ell)}$ denote the random variable defined as the maximum value over ℓ independent drawings of Z . When ℓ goes to infinity, the limiting distribution of $Z^{(\ell)}$ is the cumulative distribution $P(Z^{(\ell)} < z) \xrightarrow{\ell \rightarrow \infty} G_{\xi, \mu, \sigma}(z)$, expressed as one of the two parametric models:*

$$G_{\xi, \mu, \sigma}(z) = \exp \left\{ \begin{array}{ll} \left(1 + \xi \frac{z - \mu}{\sigma}\right)_+^{-1/\xi} & \text{if } \xi \neq 0 \\ -\exp\left(\frac{\mu - z}{\sigma}\right) & \text{otherwise} \end{array} \right\} \quad (3)$$

with $\mu \in \mathbb{R}$ a location parameter, $\sigma \in \mathbb{R}_+$ a dispersion parameter and $\xi \in \mathbb{R}$ a shape parameter referred to as extreme value index.

Overall, the EVT framework provides a general parametric model for extreme events associated with a random variable, independent of the distribution of Z itself. The universality of these models reflects the fact that modeling the extreme events associated with a distribution relies only on the behavior of its tail. This tail can take one of three forms: (i) an exponential tail ($\xi = 1$, corresponding to the Gumbel distribution); (ii) a heavy tail ($\xi > 0$, corresponding to the Fréchet distribution); or (iii) a bounded tail ($\xi < 0$, corresponding to the Weibull distribution). Despite its applicability, EVT has, to the best of our knowledge, seen limited use in machine learning, with notable exceptions in the area of anomaly detection [29, 30, 31].

3 SPADE Overview

Aimed at OOD detection, SPADE proposes a formal characterization of OOD concerning a trained teacher model and the associated latent representation on the one hand and the training distribution (hereafter in-distribution, ID) on the other hand. This characterization relies on generalized extreme value (GEV) models, which allow for detecting and rejecting out-of-ID samples. For simplicity and by abuse of language, out-of-ID samples are referred to as OOD in the following. Interestingly, under mild assumptions, the GEV models also allow for detecting adversarial samples. The abstaining classifier, equipped with the GEV-based detection tests, provides probabilistic guarantees of OOD and adversarial sample rejection, subject to a lower bound on the magnitude of the adversarial perturbation.

3.1 EVT-based Characterization of OOD

The distance-based OOD detection literature (see e.g., [12, 14, 13]) suggests that a sample is *likely* to be an OOD sample if it is *sufficiently* distant from the training samples of all (or most) classes in the latent space.

In SPADE, this process is reformulated using generalized extreme value models, directly yielding the probability for a sample to be OOD. Let (X, Y) denote the random variable following the joint distribution $P_{X, Y}$. For $Y = c$, let Z_c be the random variable defined as the distance between $h(X)$ and its k -th nearest neighbor in latent distance, belonging to \mathcal{D} with same class c . By definition, the limiting distribution of the maxima of Z_c follows a Generalized Extreme Value model noted $G^{(c)}$, with $Pr(Z_c^{(\ell)} < v) \xrightarrow{\ell \rightarrow \infty} G^{(c)}(v)$.

For each sample \mathbf{x} in the instance space and each class c , let z_c be defined as $\|h(\mathbf{x}) - h(\mathbf{x}_{knn,c})\|$ with $\mathbf{x}_{knn,c}$ the k -th nearest neighbor of \mathbf{x} in latent space, such that $(\mathbf{x}_{knn,c}, c)$ belongs to \mathcal{D} . As the true label of \mathbf{x} is unknown at inference time, the proposed OOD test retains the lowest probability of \mathbf{x} being OOD according to all $G^{(c)}$:

Definition 4 (OOD test). *Let \mathbf{x} denote an instance in \mathcal{X} with z_c its Euclidean latent distance to its k -nearest neighbor of class c in the training set \mathcal{D} . The probability of \mathbf{x} being an OOD sample, noted $OOD(\mathbf{x})$, is defined as:*

$$OOD(\mathbf{x}) = \min_{c \in \mathcal{Y}} G^{(c)}(z_c) \quad (4)$$

In other words, \mathbf{x} is considered to be OOD if it is extreme concerning all GEV models associated with the different classes.

The decision to consider a separate GEV $G^{(c)}$ for each class c is intended to address situations where classes exhibit different levels of variation in the latent space. In such cases, considering a single GEV model for all classes could result in either erroneously rejecting samples from a class with high variation or incorrectly accepting OOD samples from a class with low variation.

Since the OOD test depends on the classifier's latent representation, one might wonder to what extent different tests based on different classifiers are consistent (as discussed further in Section 5).

3.2 Abstaining Classifier on OOD Samples

A classifier abstaining on OOD samples is built as follows.

Definition 5 (Abstaining classifier). *Given teacher f and confidence $1 - \tau$, with $0 < \tau < 1$, classifier f_τ abstains from making predictions on a sample \mathbf{x} if \mathbf{x} is considered to be extreme with probability at least $1 - \tau$ w.r.t. its candidate class $c = f(\mathbf{x})$. With same notations as above:*

$$f_\tau(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{if } z_c \leq G^{(c)^{-1}}(1 - \tau) \\ \text{abstain} & \text{otherwise} \end{cases} \quad (5)$$

where z_c is the distance between $h(\mathbf{x})$ and its nearest neighbor of class c in \mathcal{D} .

The abstention test embedded in f_τ is more precise than the OOD test (Def. 4) because it incorporates the additional information of $f(\mathbf{x})$. Note, however, that both tests coincide under the common assumption that an OOD sample is closer to the samples of its own class, all else being equal.

3.3 Abstaining classifier with provable guarantees w.r.t. adversarial examples

Let us consider the adversarial example \mathbf{x} built by perturbing a training sample noted \mathbf{x}^* of class c , and let $c' = f(\mathbf{x}) \neq c$ be its (wrong) class according to f . Under mild assumptions, it is shown that the abstaining classifier f_τ abstains on adversarial sample \mathbf{x} with probability $1 - \tau$, subject to a lower bound on its perturbation amplitude.

Let $G^{(c,c')}$ denote the generalized extreme value model associated with the *minimum* latent distance among pairs of examples $(\mathbf{x}, \mathbf{x}')$ respectively belonging to class c and c' :

$$G^{(c,c')}(v) = Pr\left(\|h(X) - h(X')\| > v \mid (X, Y) \sim P_{X,Y}, (X', Y') \sim P_{X,Y}, Y = c, Y' = c'\right)$$

Theorem 1. *Let us assume that the latent embedding h is K -Lipschitz. Let \mathbf{x} be an adversarial sample built by perturbation of a training sample \mathbf{x}^* of class c , with perturbation amplitude ε ($\|\mathbf{x} - \mathbf{x}^*\| < \varepsilon$), and let $f(\mathbf{x}) = c' \neq c$. Let \mathbf{x}'^* of class c' denote the k -th nearest training sample in \mathcal{D} of \mathbf{x} . Then, with probability at least $1 - \tau$ either f_τ abstains on \mathbf{x} , or perturbation ε is greater than the following lower bound:*

$$\varepsilon \geq \frac{1}{K} \left(G^{(c,c')}^{-1}(1 - \tau) - G^{(c)}^{-1}(1 - \tau) \right) \quad (6)$$

Proof. If $\|h(\mathbf{x}) - h(\mathbf{x}^*)\| > (G^{(c)})^{-1}(1 - \tau)$, then f_τ abstains on \mathbf{x} (Def. 5). Otherwise,

$$\begin{aligned} \|h(\mathbf{x}^*) - h(\mathbf{x}'^*)\| &\leq \|h(\mathbf{x}^*) - h(\mathbf{x})\| + \|h(\mathbf{x}) - h(\mathbf{x}'^*)\| \\ &\leq K\varepsilon + (G^{(c)})^{-1}(1 - \tau) \end{aligned} \quad (7)$$

Moreover, with probability $1 - \tau$,

$$(G^{(c,c')})^{-1}(1 - \tau) \leq \|h(\mathbf{x}^*) - h(\mathbf{x}'^*)\| \quad (8)$$

Putting together Eqs. 7 and 8 concludes the proof. \square

3.4 Estimating the GEV Models

The estimation of the GEV models involved in SPADE is detailed in the case of $G^{(c)}$, defining its approximation $\widehat{G}^{(c)}$. The same process is used to learn an approximation of the $G^{(c,c')}$ models, with the only difference being that the considered extreme values are minima instead of maxima.

After [30], a straightforward approximation of an EVD proceeds by sampling the extreme events along the block maxima method. This fitting process is, however, found to be sensitive to the number of blocks and the block size. The proposed approach thus is the *Peak Over Threshold* (POT) method, which relies on the Pickands-Balkema-de-Haan theorem, often referred to as the *second theorem of EVT* [32, 33]. Formally, the POT considers the occurrences of events bypassing a threshold t , noting that the distribution of these occurrences follows a Generalized Pareto Distribution (GPD):

$$F_{\xi,\mu,\sigma}(z) = Pr(Z - t > z | Z > t) = \begin{cases} 1 - \exp\left(-\frac{z-\mu}{\sigma}\right) & \text{if } \xi = 0 \\ 1 - \left(1 + \frac{\xi(z-\mu)}{\sigma}\right)^{-1/\xi} & \text{otherwise} \end{cases} \quad (9)$$

Informally, for each class c , POT proceeds by fitting a GPD model to samples z_c that are over a threshold t_c . Formally, for each sample (\mathbf{x}_i, c) in \mathcal{D} , let $z_i \in \mathbb{R}$ be the distance of \mathbf{x}_i to its k -th nearest neighbor belonging to class c , in latent space. Let \mathcal{D}_c be the set including all such z_i for $z_i > t_c$.

The parameters (μ_c, σ_c, ξ_c) of model $\widehat{G}^{(c)}$ are learned by maximum likelihood estimation (MLE) on \mathcal{D}_c :

$$(\mu_c, \sigma_c, \xi_c) = \arg \max_{\mu, \sigma, \xi} \sum_{z \in \mathcal{D}_c} \mathcal{L}_{\mu, \sigma, \xi}(z) \quad (10)$$

with \mathcal{L} the log-probability density function of the GPD. MLE is preferred over alternative methods, e.g. the method of moments, due to its higher robustness and efficiency [30].

Eventually, model $\widehat{G}^{(c)}$ approximately characterizes whether a given sample is OOD with respect to class c with confidence $1 - \tau$. The OOD test (Def. 4) is accordingly approximated as:

$$\widehat{OOD}(\mathbf{x}) = \min_{c \in \mathcal{Y}} \widehat{G}^{(c)}(z_c) \quad (11)$$

3.5 Discussion

SPADE retains the main benefits of distance-based OOD detection approaches, being agnostic to the structure of the OOD distribution and easy to implement. Furthermore, being grounded in the EVT first principles, it enables estimating the probability for a sample to be OOD. Lastly, the abstaining classifier f_τ can also reject adversarial samples, subject to a lower bound on the magnitude of adversarial perturbations.

A potential weakness is the complexity of approximating GEV models, which is quadratic with respect to the number n of samples. This raises the question of whether stable and accurate GEV approximations can be achieved when aggressively subsampling the training set. A second issue pertains to the effectiveness of the lower bound used in rejecting adversarial samples (Eq. 6); specifically the question is whether $G^{(c,c')^{-1}}(1 - \tau) - G^{(c)^{-1}}(1 - \tau)$ is strictly positive for non-trivial confidence levels $1 - \tau$.² The SPADE method can be summarized in the algorithm below:

²Note that the requirement for $(G^{(c,c')^{-1}}(1 - \tau) - G^{(c)^{-1}}(1 - \tau))$ to be sufficiently large is reminiscent of the clustering assumption that underpins semi-supervised learning [34].

Algorithm 1 SPADE. Learning EVT models

```

1: Input: training set  $\mathcal{D}$ , integer  $k$ , threshold  $t > 0$ 
2: Output: GEV distributions  $(\hat{G}^{(c)})$  of each class  $c$ 
3:
4: for  $c \in \mathcal{Y}$  do
5:    $\mathcal{T}^{(c)} = \{\}$ 
6:   for  $(\mathbf{x}, c) \in \mathcal{D}$  do
7:     Compute normalized activation:  $\mathbf{z} = h(\mathbf{x}) / \|h(\mathbf{x})\|$ 
8:     Compute distance to  $k$ -th nearest neighbor in class  $c$ :  $v = \min_{(\mathbf{x}', c) \in \mathcal{D}}^{(k)} (\|\mathbf{z} - \mathbf{z}'\|)$ 
9:     if  $v > t$ :  $\mathcal{T}^{(c)} = \mathcal{T}^{(c)} \cup \{v\}$  end if
10:  end for
11:  Fit an Extreme Value Distribution  $(\hat{G}^{(c)})$  from extreme samples of  $\mathcal{T}^{(c)}$ 
12: end for
13: return  $(\hat{G}^{(c)})_{c \in \mathcal{Y}}$ 

```

4 Experimental Setting

This section outlines the experimental setup used to evaluate SPADE in comparison to the state of the art. All experiments were conducted on Tesla A100 80GB GPUs. Further details are provided in the supplementary material (SM).

Goals. The experiments aim to empirically address four questions: the performance of the OOD detection tests based on the GEV models, particularly in comparison with distance-based approaches (Q1); their sensitivity with respect to the considered teacher (Q2); similarly, the performance of the adversarial sample detection test based on the GEV models, compared with state-of-the-art methods (Q3); the computational complexity and stability of the approximate GEV models embedded in SPADE (Q4).

The performance of SPADE is compared to five established baselines: the seminal MSP [4], ODIN [5], MDS [35], KNN [13], and its extension CIDER [15] (further discussed in Section 6).

Metrics. The comparative evaluation is conducted using the OpenOOD framework [36], with performance assessed by standard indicators: The *Area Under the ROC Curve (AUC)* measures the average rate of correct OOD/adversarial sample detection across all confidence levels $2 - \tau$, corresponding to the true positive rate, as described in [36]. The *FPR95 indicator* represents the fraction of true samples misclassified as OOD (respectively, adversarial) when the detection threshold ensures 95% of OOD (resp., adversarial) samples are correctly rejected.

Benchmarks. Three medium- and large-sized datasets are considered: CIFAR-10, CIFAR-100 [37] and ImageNet-1K (using the ILSVRC2012 version). Three types of neural architectures are used to assess the generality of the SPADE approach: ResNet [19], ViT [20], and VGG [38]. The distance in latent space between a sample and its k -th nearest neighbor is calculated using the normalized L_2 distance, following [13]. OOD samples are sourced from near-OOD and far-OOD datasets, following [22, 25, 23, 24, 26], as detailed in the supplementary material (SM). Adversarial samples are generated by perturbing training samples using FGSM [17] attacks, with a perturbation amplitude ε ranging from 0.001 to 0.004.

5 Experimental Results

OOD Detection (Q1). The performance of SPADE is illustrated in Table 1, focusing on the representative case of ImageNet-1K and considering a ResNet teacher. For the considered near-OOD datasets, the best method is MSP, whereas for far-OOD datasets, the best method is KNN. In all cases but one, SPADE-ResNet is slightly outperformed by KNN. In terms of rank (determined by the average of AUC and FPR95), SPADE ranks second best on both near- and far-OOD datasets.

Sensitivity of OOD Detection (Q2). The impact of the considered teacher on the performance of the SPADE OOD test is illustrated in Table 2, comparing SPADE built on teachers ResNet, VGG, ViT-B16 with the CIDER baseline on CIFAR-10. These results show that the detection accuracy indeed depends on the teacher

Table 1: OOD detection on ImageNet-1K: performance of SPADE (with ResNet teacher) compared to that of baselines MSP, ODIN, MDS and KNN in terms of AUC (the greater the better) and FPR95 (the lower the better; best performances in bold). The rank, averaged over far and near OOD datasets, is computed after the half sum of AUC and FPR95.

	Near OOD				Far OOD						Rank
	SSB Hard		NINCO		iNaturalist		Textures		OpenImages-O		
	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	
MSP [1]	72.53	74.43	80.66	57.72	87.78	44.08	82.81	59.16	85.21	49.62	3
ODIN [2]	72.51	77.36	77.55	70.83	89.51	41.46	87.02	56.58	86.33	54.10	4
MDS [3]	52.15	90.46	68.49	71.66	76.49	56.07	94.11	27.07	77.68	59.66	5
KNN [4]	62.80	84.08	79.30	58.92	84.62	42.39	96.06	23.39	86.38	44.24	1
SPADE	61.91	85.27	77.99	61.04	85.26	44.84	95.86	24.63	85.79	46.33	2

Table 2: OOD detection on CIFAR-10: performance of SPADE with teachers ResNet, VGG and ViT-B16, compared to that of baseline CIDER.

	TIN		MNIST		SVHN		Textures		Places-365	
	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow
CIDER ($d=512$)	71.56	70.34	68.84	71.86	57.51	78.43	71.06	70.70	71.73	69.97
ResNet-18 ($d=512$)	90.41	32.14	93.646	21.91	92.17	22.94	91.97	25.68	91.03	30.41
VGG-16 ($d=512$)	76.04	66.84	89.37	35.12	81.53	44.06	80.47	50.73	72.14	77.27
ViT-B16 ($d=384$)	96.27	20.18	94.52	12.23	81.78	34.83	99.97	00.06	99.67	00.51

and its latent space, with SPADE-ViT-B16 dominating ResNet (except on SVHN) and ResNet strongly dominating VGG. Still, the performance does not merely depend on the size of the latent space. The discrepancy between the AUC and FPR95 indicators suggests that the optimal detection threshold varies depending on the dataset. Notably, all SPADE OOD detection tests outperform CIDER, which is based on KNN and specifically targets OOD detection. We shall come back to this remark in Section 6.

Detection of Adversarial Samples (Q3). Table 3 reports the performance of SPADE (with a ResNet teacher) on the representative cases of CIFAR-10 and CIFAR-100, compared with MSP, MDS, KNN and CIDER. For all perturbation amplitudes, SPADE ranks first w.r.t. AUC (and second w.r.t. FPR95). In terms of FPR95, KNN ranks first on CIFAR-10, while MSP ranks first on CIFAR-100. Overall, SPADE behaves on par with, or better than, OOD detection methods w.r.t. the detection of adversarial examples. The slight AUC improvement suggests that SPADE may capture more subtle differences between in-distribution and adversarial samples. Conversely, the high FPR95 values suggest that SPADE tends to be overly cautious, rejecting true samples at the level of confidence where 95% adversarial samples are rejected.

Computational Frugality: Stability of GEV Models and SPADE Performances (Q4). The stability of the GEV models with respect to the fraction of training samples used to estimate the $\hat{G}^{(c)}$ hyperparameters is illustrated in Fig. 1, using the ResNet teacher’s latent space on CIFAR-100. The figure shows: i) a low sensitivity of the tail index (Fig. 1.(a)); ii) a decrease in μ with the sampling rate (Fig. 1.(b)); iii) a low sensitivity of the dispersion parameter σ (Fig. 1.(c)). Complementary results are provided in the SM, confirming that the lower bound on the adversarial perturbation amplitude is effectively positive.

As expected, the stability of the $\hat{G}^{(c)}$ models results in stable OOD detection performances (AUC and FPR95) when subsampling the training set. Quite the contrary, the OOD detection performances of KNN are significantly deteriorated when aggressively subsampling the training set, particularly so on near-OOD (Fig. 2).

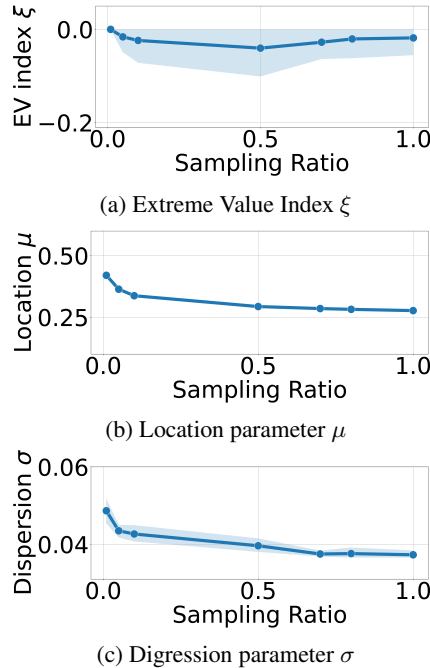


Figure 1: Stability of EVT parameter estimation wrt sampling ratio and estimation variance for one class of CIFAR-100 on ResNet-18.

Table 3: Adversarial samples detection on CIFAR-10 and CIFAR-100, with perturbation amplitude from .001 to .004: comparison of SPADE (ResNet teacher) with baselines MSP, ODIN, MDS and KNN in terms of AUC (the greater the better) and FPR95 (the lower the better; best in bold).

		$\epsilon = 0.001$		$\epsilon = 0.002$		$\epsilon = 0.003$		$\epsilon = 0.004$		Average	
		AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow	AUC \uparrow	FPR95 \downarrow
CIFAR-10	MSP [1]	81.68	79.08	81.90	78.70	82.06	79.14	82.20	78.10	81.96	78.76
	MDS [3]	81.46	69.34	81.51	68.76	81.57	69.04	81.61	69.42	81.54	69.14
	KNN [4]	85.65	54.46	85.75	54.19	85.85	53.91	85.92	54.29	85.79	54.21
	CIDER [5]	85.46	55.68	85.50	54.79	85.53	54.93	85.60	54.57	85.52	54.99
	SPADE	85.96	55.02	86.06	54.51	86.15	54.40	86.24	53.78	86.10	54.43
CIFAR-100	MSP [1]	83.24	51.84	83.39	50.57	83.54	49.94	83.68	49.64	83.46	50.50
	MDS [3]	60.34	82.93	60.29	83.01	60.25	82.60	60.19	83.01	60.27	82.89
	KNN [4]	83.54	56.83	83.67	56.39	83.79	55.64	83.89	55.13	83.72	56.00
	CIDER [5]	82.75	63.17	82.85	63.57	82.95	63.17	83.06	62.34	82.90	63.06
	SPADE	84.33	53.13	84.45	52.84	84.56	52.44	84.66	52.44	84.50	52.72

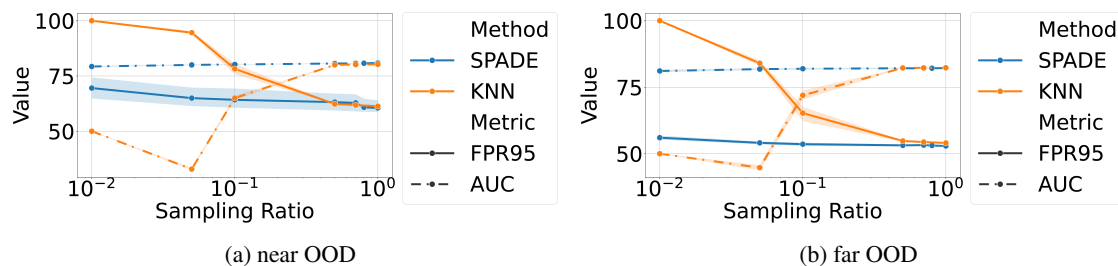


Figure 2: Sensitivity analysis of the OOD detection on CIFAR-100 w.r.t. the subsampling rate of the training set: AUC (dashed line) and FPR95 (solid line) performances for SPADE (in blue) and KNN [13] (in orange; better seen in color).

6 Position w.r.t. related work

The robustness lack of machine learning models with respect to adversarial and OOD examples is widely recognized as a major obstacle for ML applications in safety-critical domains [1, 39].

With respect to OOD examples, their detection takes inspiration from several areas of ML, ranging from learning with rejection [40], to anomaly detection [41], novelty detection [42, 43], and open set recognition [44]. To the best of our knowledge, the OOD detection problem was first formalized by [45]; the early MSP approach, observing the margin between the logits of the trained teacher and exploiting the fact that it behaves differently for in-distribution and OOD samples [4], still is among the most effective ones (Table 1). Along the same line, ODIN exploits the gradient information to separate in- and out-of-distribution samples [5].

Some approaches address OOD detection as yet another supervised learning issue, treating OOD samples as belonging to an additional class and utilizing them to (re-)train the model [46, 47] (see also [6]).

Quite a few other methods are based on the so-called manifold assumption; the challenge, then, is to identify the representation that best characterizes the manifold on which real samples lie. One option is to consider the latent representation of an auto-encoder (AE) trained solely on real samples. As empirically shown by e.g. [9, 10, 46], the reconstruction error of OOD samples through the AE is much higher than for real samples; an effective OOD test can thus be based on this error. On the positive side, this reconstruction error defines a general criterion and does not depend on the teacher; on the negative side, it does not leverage class information. Another option is based on modeling the behavior of true samples in the penultimate layer of the neural net, using probabilistic models [11, 10, 9].

Finally, the option most closely related to SPADE is to consider the latent representation of the teacher itself, as is done in distance-based OOD detection approaches [12, 35], particularly in KNN [13]. The difference is that KNN exploits the distance z of a sample to its k -th nearest neighbor in the training set, while SPADE exploits $\hat{G}^{(c)}(z_c)$. A tentative interpretation for the better performance of KNN (Table 1) is the superior bias-variance trade-off in the empirical test based on z compared to the test based on $\hat{G}^{(c)}(z_c)$. While learning a parametric GEV model achieves

some regularization, this model is only trained from examples within the class c , meaning it operates with one or two orders of magnitude less data.

The problem of dealing with adversarial examples differs from that of detecting OOD examples, as adversarial examples are deliberately crafted to deceive the teacher [48, 17, 49, 50]. Knowing their structure allows for the design of specific defense strategies, such as adversarial training [17, 49], which incorporates adversarial examples into the training process [51, 52, 53]. Other notable defense strategies include adversarial architectures [54, 55], adversarial regularization [56, 57], and data augmentation methods [58, 59].

In contrast, SPADE neither requires additional information nor retrains the classifier to defend against adversarial examples. It employs the same agnostic strategy for both adversarial and OOD threats: it characterizes the true samples and abstains from making a decision if the example in question appears extreme compared to the true (ID) samples, based on the chosen confidence level.

7 Conclusion and Perspectives

The main contribution of the paper, SPADE (*Sample-efficient Probabilistic Detection using Extreme value theory*), is a formal test designed to detect examples appearing to be extreme w.r.t. training samples, enabling the classifier to abstain from making decisions on such extreme examples. This test’s ability to accurately detect OOD and adversarial samples has been empirically demonstrated, with similar performances as the prominent state of the art approaches. Furthermore, the computational complexity and the stability of the proposed detection test have been empirically established.

The proposed test, similar to distance-based OOD detection approaches, exploits the latent distance between the given example and its nearest neighbor in the training set. Its originality lies in leveraging Extreme Value Theory [18] to provide a formal characterization of the training samples. This characterization offers two key benefits: first, it yields provable guarantees for detecting adversarial examples, subject to the adversarial perturbation amplitude to be lower bounded; second, it provides some new hints into the key aspects governing the teacher robustness.³

This approach opens several perspectives, related with making classifiers more robust and better understanding the key robustness factors. A short term perspective is to extend the generalized extreme value test to some of the empirical criteria used in the OOD detection literature; one such criterion is the score margin involved in MSP [4]. Another perspective is to enhance the classifier training loss to favor the robustness of the latent space, e.g. to consider the optimization of the Lipschitz constant of the classifier embedding besides its variation and informativeness as done in CIDER [15].

Our long-term goal is to investigate whether *safe* example behaviors can be identified in the latent space and whether these behaviors can be certified, as a step toward the certification of neural networks.

References

- [1] Mohammadreza Salehi, Hossein Mirzaei, Dan Hendrycks, Yixuan Li, Mohammad Hossein Rohban, and Mohammad Sabokrou. A unified survey on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future challenges. *TMLR*, 2022.
- [2] Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyun Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan Hendrycks, Yixuan Li, and Ziwei Liu. Openood: Benchmarking generalized out-of-distribution detection. *NeurIPS*, 2022.
- [3] Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, and Liwei Wang. Towards a theoretical framework of out-of-distribution generalization. *NeurIPS*, 2021.
- [4] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *ICLR*, 2017. [Alias \(1\) used in Table 1.](#)
- [5] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *ICLR*, 2018. [Alias \(2\) used in Table 1.](#)
- [6] Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. *CVPR*, 2020.

³For instance, CIDER [15] involves the optimization of the variation and informativeness of the teacher latent space. The SPADE analysis suggests that besides these two factors, the regularity of the teacher (its Lipschitz constant) also matters. A possible interpretation for why SPADE outperforms CIDER is that the optimization of the variation and informativeness might adversely affect this Lipschitz constant.

- [7] Wenyu Jiang, Yuxin Ge, Hao Cheng, Mingcai Chen, Shuai Feng, and Chongjun Wang. Read: Aggregating reconstruction error into out-of-distribution detection. *AAAI*, 2023.
- [8] Jingyao Li, Pengguang Chen, Shaozuo Yu, Zexin He, Shu Liu, and Jiaya Jia. Rethinking out-of-distribution (ood) detection: Masked image modeling is all you need, 2023.
- [9] Jie Ren, Peter J. Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. *NeurIPS*, 2019.
- [10] Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection. *NeurIPS*, 2021.
- [11] Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don’t know by virtual outlier synthesis. *ICLR*, 2022.
- [12] Nicolas Papernot and Patrick D. McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *ArXiv*, 2018.
- [13] Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. *ICML*, 2022. [Alias \(4\) used in Table 1.](#)
- [14] Adam Dziedzic, Stephan Rabanser, Mohammad Yaghini, Armin Ale, Murat A. Erdogdu, and Nicolas Papernot. p-dknn: Out-of-distribution detection through statistical testing of deep representations. *ArXiv*, 2022.
- [15] Yifei Ming, Yiyu Sun, Ousmane Dia, and Yixuan Li. How to exploit hyperspherical embeddings for out-of-distribution detection? *ICLR*, 2023. [Alias \(5\) used in Table 1.](#)
- [16] Badr Youbi Idrissi, Diane Bouchacourt, Randall Balestrieri, Ivan Evtimov, Caner Hazirbas, Nicolas Ballas, Pascal Vincent, Michal Drozdal, David Lopez-Paz, and Mark Ibrahim. Imagenet-x: Understanding model mistakes with factor of variation annotations. *ICLR*, 2023.
- [17] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *ICLR*, 2015.
- [18] Ronald Aylmer Fisher and L. H. C. Tippett. Limiting forms of the frequency distribution of the largest or smallest member of a sample. *Mathematical Proceedings of the Cambridge Philosophical Society*, 1928.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, 2016.
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
- [21] Sebastian Farquhar and Yarin Gal. What ‘out-of-distribution’ is and is not. *NeurIPS ML Safety Workshop*, 2022.
- [22] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. *CVPR*, 2014.
- [23] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. *CVPR*, 2018.
- [24] Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching. *CVPR*, 2022.
- [25] Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: a good closed-set classifier is all you need? *ICLR*, 2022.
- [26] Julian Bitterwolf, Maximilian Mueller, and Matthias Hein. In or out? fixing imagenet out-of-distribution detection evaluation. *ICML*, 2023.
- [27] Laurens Haan and Ana Ferreira. *Extreme Value Theory: An Introduction*. Springer Series in Operations Research and Financial Engineering. Springer, New York, 2006.
- [28] Boris Gnedenko. Sur la distribution limite du terme maximum d’une serie aleatoire. *Annals of mathematics*, 1943.
- [29] Mark Smith, Steven Reece, Stephen J. Roberts, and Iead Rezek. Online maritime abnormality detection using gaussian processes and extreme value theory. *ICDM*, 2012.
- [30] Alban Siffer, Pierre-Alain Fouque, Alexandre Termier, and Christine Largouet. Anomaly detection in streams with extreme value theory. *KDD*, 2017.
- [31] Joshua P. French, Piotr Kokoszka, Stilian A. Stoev, and Lauren Hall. Quantifying the risk of heat waves using extreme value theory and spatio-temporal functional data. *Comput. Stat. Data Anal.*, 2019.

- [32] A. A. Balkema and L. de Haan. Residual Life Time at Great Age. *The Annals of Probability*, 1974.
- [33] J. Pickands. Statistical Inference Using Extreme Order Statistics. *The Annals of Statistics*, 1975.
- [34] Philippe Rigollet. Generalization error bounds in semi-supervised classification under the cluster assumption. *JMLR*, 2007.
- [35] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *NeurIPS*, 2018. [Alias \(3\) used in Table 1.](#)
- [36] Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyou Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Yixuan Li, Ziwei Liu, Yiran Chen, and Hai Li. Openood v1.5: Enhanced benchmark for out-of-distribution detection. *arXiv*, 2023.
- [37] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technical Report*, 2009.
- [38] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *ACPR*, 2015.
- [39] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *IJCV*, 2024.
- [40] Peter L. Bartlett and Marten H. Wegkamp. Classification with a reject option using a hinge loss. *JMLR*, 2008.
- [41] Saikiran Bulusu, Bhavya Kailkhura, Bo Li, Pramod K. Varshney, and Dawn Song. Anomalous example detection in deep learning: A survey. *IEEE Access*, 2020.
- [42] Markos Markou and Sameer Singh. Novelty detection: a review—part 1: statistical approaches. *Signal Processing*, 2003.
- [43] Markos Markou and Sameer Singh. Novelty detection: a review—part 2: neural network based approaches. *Signal Processing*, 2003.
- [44] T. E. Boulton, S. Cruz, A.R. Dhamija, M. Gunther, J. Henrydoss, and W.J. Scheirer. Learning and the unknown: Surveying steps toward open world recognition. *AAAI*, 2019.
- [45] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *CVPR*, 2015.
- [46] Xuefeng Du, Xin Wang, Gabriel Gozum, and Yixuan Li. Unknown-aware object detection: Learning what you don’t know from videos in the wild. *CVPR*, 2022.
- [47] Jingyang Zhang, Nathan Inkawhich, Randolph Linderman, Yiran Chen, and Hai Li. Mixture outlier exposure: Towards out-of-distribution detection in fine-grained environments. *WACV*, 2022.
- [48] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *ICLR*, 2014.
- [49] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *ICLR*, 2018.
- [50] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. *ICML*, 2020.
- [51] Ziang Yan, Yiwen Guo, and Changshui Zhang. Deep defense: Training dnns with improved adversarial robustness. *NeurIPS*, 2018.
- [52] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. *ICML*, 2019.
- [53] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial robustness requires revisiting misclassified examples. *ICLR*, 2020.
- [54] Ramtin Hosseini, Xingyi Yang, and Pengtao Xie. Dsrna: Differentiable search of robust neural architectures. *CVPR*, 2020.
- [55] Shihua Huang, Zhichao Lu, Kalyanmoy Deb, and Vishnu Naresh Boddeti. Revisiting residual networks for adversarial robustness. *CVPR*, 2023.
- [56] Chengzhi Mao, Ziyuan Zhong, Junfeng Yang, Carl Vondrick, and Baishakhi Ray. Metric learning for adversarial robustness. *NeurIPS*, 2019.
- [57] Hong Liu, Zhun Zhong, Nicu Sebe, and Shin’ichi Satoh. Mitigating robust overfitting via self-residual-calibration regularization. *Artificial Intelligence*, 2023.

- [58] Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled data improves adversarial robustness. *NeurIPS*, 2019.
- [59] Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion models further improve adversarial training. *ICML*, 2023.