

Non-Parametric Probabilistic Robustness: A Conservative Metric with Optimized Perturbation Distributions

Zheng Wang, Yi Zhang, Siddartha Khastgir, Carsten Maple, Xingyu Zhao*
WMG, University of Warwick, Coventry CV4 7AL, United Kingdom

{Zheng.Wang.3,Yi.Zhang.16,S.Khastgir.1,CM,xingyu.zhao}@warwick.ac.uk

Abstract

Deep learning (DL) models, despite their remarkable success, remain vulnerable to small input perturbations that can cause erroneous outputs, motivating the recent proposal of probabilistic robustness (PR) as a complementary alternative to adversarial robustness (AR). However, existing PR formulations assume a fixed and known perturbation distribution, an unrealistic expectation in practice. To address this limitation, we propose non-parametric probabilistic robustness (NPPR), a more practical PR metric that does not rely on any predefined perturbation distribution. Following the non-parametric paradigm in statistical modeling, NPPR learns an optimized perturbation distribution directly from data, enabling conservative PR evaluation under distributional uncertainty. We further develop an NPPR estimator based on a Gaussian Mixture Model (GMM) with Multilayer Perceptron (MLP) heads and bicubic up-sampling, covering various input-dependent and input-independent perturbation scenarios. Theoretical analyses establish the relationships among AR, PR, and NPPR. Extensive experiments on CIFAR-10, CIFAR-100, and Tiny ImageNet across ResNet18/50, WideResNet50 and VGG16 validate NPPR as a more practical robustness metric, showing up to 40% more conservative (lower) PR estimates compared to assuming those common perturbation distributions used in state-of-the-arts.

1. Introduction

Deep learning (DL) models, despite their success across perception, language, and control tasks, are known to be vulnerable to small input perturbations that can drastically alter their outputs. Such perturbations, commonly referred to as adversarial examples (AEs), reveal a fundamental weakness in modern DL models. Since the first discovery of the phenomenon of AEs [11, 21, 24], robustness has emerged as one of the most actively studied properties.

*Corresponding to: X. Zhao, xingyu.zhao@warwick.ac.uk

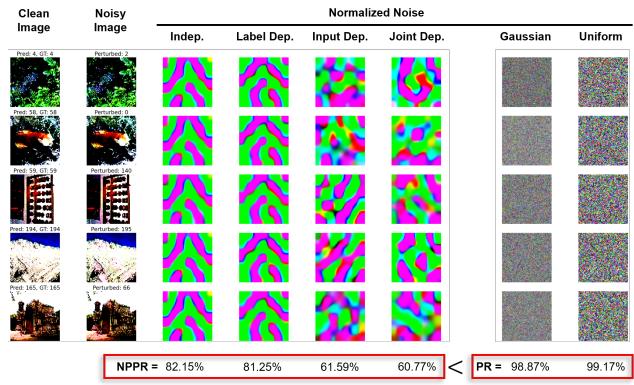


Figure 1. **Examples of perturbation results.** We visualize perturbations (normalized to the range (0, 1) for visibility) generated by our optimisation pipeline using ResNet-18 on TinyImageNet under four different dependency settings. As shown, the optimized perturbations differ substantially from those produced by Gaussian or uniform noise (that are commonly assumed as fixed perturbation distributions in state-of-the-arts). Our proposed NPPR metric yields more conservative (i.e., lower) estimates on PR.

Among them, adversarial robustness (AR) [4, 13, 20, 40] is arguably the most extensively studied, focusing on a model’s ability to withstand *deterministic* and *worst-case*¹ AEs deliberately crafted by malicious attackers from a security perspective. In contrast, probabilistic robustness (PR) [1, 7, 22, 23, 25, 26, 28, 31–33, 35, 36, 38, 39] has recently been proposed as a complementary notion, not only from the reliability perspective (concerning average risk [28, 31, 39]) but also as an extension of the security viewpoint, aiming to quantify the *likelihood* that a DL model maintains correct behavior under *random* perturbations. Such random perturbations may arise from natural stochastic sources, e.g., sensor white noise in benign operational environments, or from unsophisticated attackers who rely on brute-force random-noise strategies rather than care-

¹The terms “worst-case” and “deterministic” indicate that AR studies typically aim to find the optimized AE, e.g., the one that maximizes loss or lies closest to the original input within a norm-ball.

fully optimized manipulations (i.e., attacks typically proposed in AR studies).

Intuitively, PR addresses the question: “What is the likelihood of encountering AEs under stochastic perturbations?”. Such stochastic perturbations should be generated from *a given distribution*, referred to as the “*input model*” in, e.g., [31, 32, 35]. Indeed, this perturbation distribution is assumed to be known and fixed in the very first formal definition of PR [31] and has since been adopted by *all* subsequent studies in both PR assessment [1, 7, 22, 25, 26, 28, 31–33, 35] and PR training [23, 28, 34, 37, 38] works. However, in practice, *the perturbation distribution is rarely known a priori*. Although Gaussian or Uniform distributions are commonly used in the literature, these choices are merely illustrative and not intended to represent any universal perturbation distributions. As noted in those PR works, the assessor is expected to specify an appropriate perturbation distribution on a case-by-case basis, supported by justified evidence from the target application—an expectation that is often infeasible in practice.

To bridge the gap arising from the unknown perturbation distribution, we introduce non-parametric probabilistic robustness (NPPR), a more practical PR metric that *does not rely on any predefined perturbation distribution*. The term “non-parametric” follows its usage in statistical modelling, referring to approaches that do not assume a fixed parametric form for the underlying distribution but instead infer it from data. Accordingly, our proposed NPPR introduces an *optimized* perturbation distribution learned from data. This enables the PR to be evaluated *conservatively*, i.e., yielding the lowest PR estimate within the admissible perturbation distribution space (see Fig. 2).

Specifically, after formally defining NPPR, we develop an NPPR estimator that fits a Gaussian Mixture Model (GMM) with Multilayer Perceptron (MLP) heads and bicubic up-sampling to optimize the perturbation distribution from data for the most conservative PR estimates. The estimator considers four dependency scenarios between perturbations, inputs, and labels. We further provide a theoretical analysis of the relationships among AR, PR, and NPPR under these scenarios. Extensive experiments on ResNet18/50, WideResNet50 and VGG16 across CIFAR-10, CIFAR-100, and Tiny ImageNet datasets are conducted to validate our approach.

In summary, our main contributions are as follows:

1. **Formal metric:** We extend the PR concept to a non-parametric setting and formally define the NPPR metric, with theoretical analysis of its relationship to AR & PR.
2. **Estimator:** We develop an NPPR estimator that accommodates various input-dependent and input-independent perturbation scenarios.
3. **Open-source repository:** All experimental details are

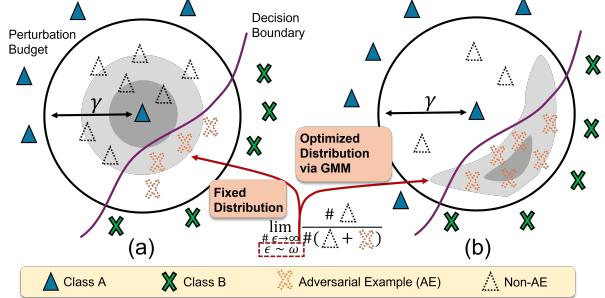


Figure 2. **Illustration of PR vs. NPPR.** Panel (a) illustrates PR, which measures the relative proportion of non-adversarial examples (Non-AEs) under a predefined fixed distribution. Panel (b) depicts NPPR, which evaluates the same metric under an optimized distribution learned via a GMM, resulting in a higher proportion of AEs and thus a more conservative robustness estimate.

provided in the supplementary materials, and a public repository will be released after the review process.

2. Preliminary and Related Works

Consider a standard classification task, where each input-label pair $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is drawn from an unknown data distribution D over $\mathcal{X} \times \mathcal{Y}$. Here, $\mathcal{X} \subseteq \mathbb{R}^d$ denotes the input space and $\mathcal{Y} = \{1, 2, \dots, C\}$ the label space. A hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$, with $h \in \mathcal{H}$, represents the classifier under study. The training set $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ consists of N i.i.d. samples drawn from D , and $\ell : \mathcal{H} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ denotes the loss function.

We denote perturbations by $\epsilon \sim \omega$, where ω is the underlying perturbation distribution, and let \mathcal{B} denote the admissible perturbation budget (e.g., an L_p ball). For convenience in our theoretical analysis, we adopt the notation \mathfrak{S} from prior work on adversarial robustness [8, 30] to represent our local robustness metric, and \mathcal{G} for corresponding global robustness metric.

2.1. Adversarial vs. Probabilistic Robustness

Although definitions of robustness vary across different DL tasks and model types, it generally refers to a model’s ability to maintain consistent predictions under small input perturbations. Typically, robustness is defined such that all inputs within a perturbation budget \mathcal{B} yield the same prediction, where \mathcal{B} usually denotes an L_p -norm ball of radius γ around an input \mathbf{x} . A perturbed input \mathbf{x}' (e.g., obtained by adding noise to \mathbf{x}) is considered an AE if its predicted label differs from that of \mathbf{x} . Evaluation methods for robustness are typically based on metrics defined over AEs [13, 40]. The formal definitions of AR and PR metrics are introduced in Def. 1 and Def. 2, respectively.

Definition 1 (Adversarial Robustness) *Given a classifier $h \in \mathcal{H}$, and the loss function ℓ , AR around the input \mathbf{x} is*

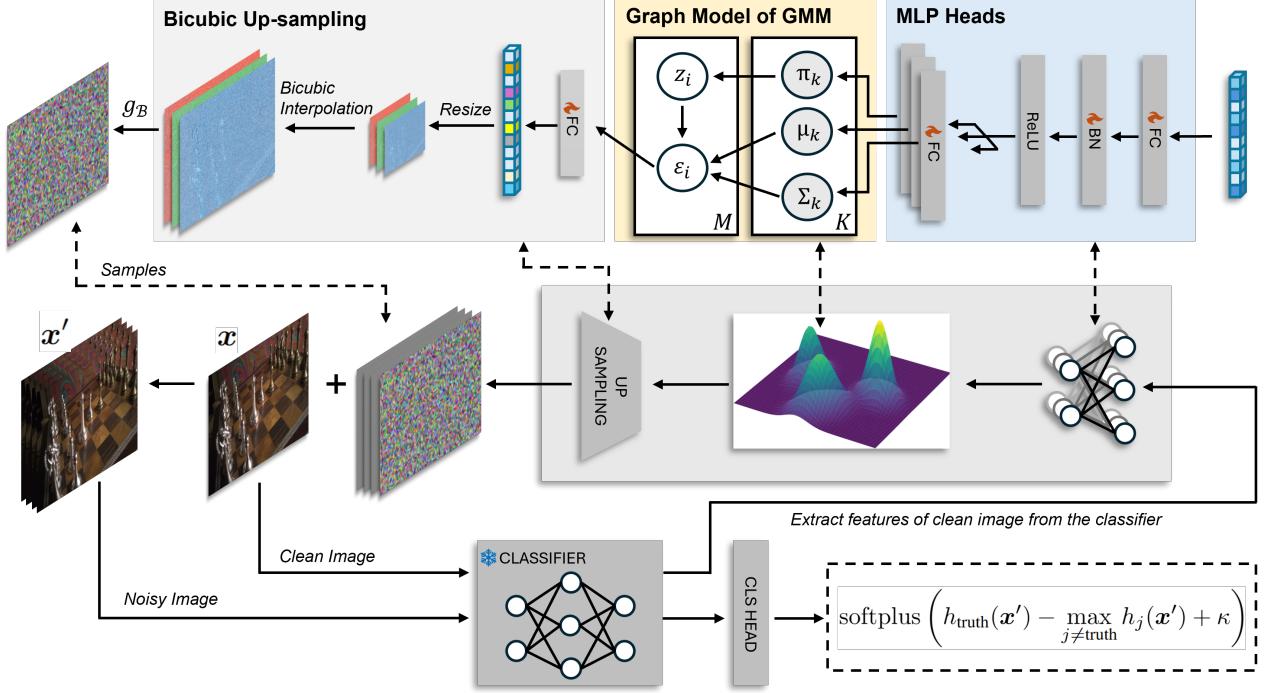


Figure 3. Training pipeline of NPPR estimator. Our training pipeline comprises three main components: (i) *MLP heads* that model the dependency structure of perturbations, (ii) a GMM for sampling latent perturbations, and (iii) *bicubic up-sampling* to map perturbations to the input space. Given a clean image \mathbf{x} , we extract intermediate features from the classifier and pass them through the MLP heads to parameterize the GMM. Perturbations are subsequently sampled from the GMM and up-sampled to the input resolution via bicubic interpolation. g_B maps the unbounded perturbations into the perturbation budget. The classifier’s logits for the perturbed input \mathbf{x}' are then used to construct a loss function based on the logit margin between the ground-truth class and the most confident non-ground-truth class. The margin parameter κ controls the scale of the gap, following the formulation introduced in the C&W attack [3].

defined as

$$\mathfrak{S}_{\text{AR}}(\mathbf{x}, y) \triangleq \sup_{\varepsilon \in \mathcal{B}} \ell(h(\mathbf{x} + \varepsilon), y), \quad (1)$$

where $\mathcal{B} = \{\varepsilon \in \mathbb{R}^d \mid \|\varepsilon\|_p \leq \gamma\}$ denotes the perturbation budget.

Definition 2 (Probabilistic Robustness) Reusing notations in Def. 1, let $\omega(\cdot | \mathbf{x})$ denotes a perturbation distribution conditioned on \mathbf{x} , whose support lies within \mathcal{B} . Let $\mathbf{1}_{\mathcal{S}(\mathbf{x})}$ be an indicator function that equals 1 if \mathcal{S} holds and 0 otherwise. Then, the PR of an input–label pair (\mathbf{x}, y) is defined as

$$\mathfrak{S}_{\text{PR}}(\mathbf{x}, y, \omega) \triangleq \mathbb{E}_{\varepsilon \sim \omega(\cdot | \mathbf{x})} [\mathbf{1}_{h(\mathbf{x} + \varepsilon) = y}] \quad (2)$$

Remark 1 (Input-dependency of perturbations) While Def. 2 provides a general formulation of PR in which the perturbation distribution ω depends on the input \mathbf{x} , in practice such perturbations may or may not depend on the specific input. Input-dependent perturbations occur when the characteristics of the noise vary with the input itself. For instance, in computer vision, noise may increase

in darker regions of an image, motion blur may depend on the object’s velocity. In contrast, input-independent perturbations remain statistically identical across all inputs, such as additive white Gaussian noise, disturbances from constant background vibration or temperature drift that affect all camera inputs equally.

AR captures the “worst-case” scenario in which the generated perturbations yield the AE that maximized the loss or closest to the input, typically requiring carefully designed adversarial attack algorithms that often rely on access to model gradients. PR, in contrast, complements AR by estimating the *likelihood* of encountering AEs when perturbations are generated stochastically, ensuring that the risk remains below an acceptable threshold rather than being exactly zero [31–33, 35, 39]. Intuitively, Def. 2 defines PR as the probability that a model’s prediction remains unchanged under random perturbations of an input \mathbf{x} . A “frequentist” interpretation of this probability is the *limiting relative frequency* of perturbations that do not alter the predicted label over infinitely independent trials [35, 39] (see Fig. 2, panel (a)). We also note that the PR definition in the literature constrains perturbations within a L_p -norm ball by

enforcing $\|\varepsilon\|_p \leq \gamma$. Equivalently, we assume a noise distribution ω with zero probability mass outside this budget.

2.2. Probabilistic Robustness Assessment

Recent years have witnessed notable progress in PR research, particularly in terms of its assessment, resulting in the development of a variety of assessment methods across different DL tasks. The earliest study [31] to formally define and evaluate PR introduced a black-box statistical estimator based on the Multi-Level Splitting method [16], which decomposes the task of estimating the probability of a rare event into several subproblems and is thus suitable for cases where PR is very high. Later, more efficient white-box estimators were developed in [26]. Zhang *et al.* [33] further investigated PR under functional perturbations such as color shifts and geometric transformations. In addition, PR has been extended to broader applications such as explainable AI [12] and text-to-image models [35]. For a comprehensive overview of PR assessment, readers are referred to [39]. Beyond assessment, several studies [23, 28, 34, 37] have focused on developing training methods to improve PR, with their optimization strategies systematically summarized in the recent benchmark study [38].

All aforementioned state-of-the-arts rely on a shared assumption that the perturbation noise follows a predefined distribution, which is rarely known in practice. In contrast, our proposed NPPR learns an optimized distribution from data, enabling a conservative evaluation that yields the lowest PR estimate within the admissible distribution space.

3. Non-Parametric Probabilistic Robustness

As illustrated in Fig. 2, an inappropriate perturbation distribution (as in existing PR studies that predefine a fixed perturbation distribution; cf. Fig. 2 (a)) may underestimate the true risk. Therefore, an optimized perturbation distribution, adaptively learned from real data, is necessary to overcome the infeasible assumption of a predefined distribution. Accordingly, we formally define our NPPR in Def. 3, in which the learned optimized perturbation distribution allows PR to be evaluated more conservatively and accurately.

Definition 3 (Non-parametric PR) Reusing the notation from Def. 2, the NPPR of an input-label pair (\mathbf{x}, y) is defined as

$$\mathfrak{S}_{\text{NPPR}}(\mathbf{x}, y) \triangleq \inf_{\omega \in P_\varepsilon} \mathbb{E}_{\varepsilon \sim \omega(\cdot | \mathbf{x})} [\mathbf{1}_{h(\mathbf{x} + \varepsilon) = y}], \quad (3)$$

where P_ε is the distribution family with support all lying within the perturbation budget \mathcal{B} .

Remark 2 Similar to Def. 2, NPPR also accommodates both input-dependent and input-independent noise distributions (see Remark 1). Based on this property, we design

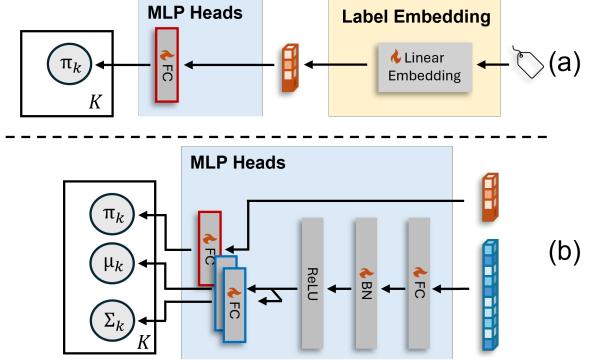


Figure 4. **Different dependency constructions.** We employ distinct MLP heads to model different dependency structures. **Panel (a)** illustrates the setting in which the perturbation distribution is conditioned solely on the ground-truth label, whereas **panel (b)** depicts the joint dependency case, where perturbations are conditioned on both the input features and labels, with the labels influencing only the mixture proportions. The label embedding in panel (b) is omitted for clarity, as it is identical to that in panel (a).

different estimators tailored to each type of dependence, enabling a more accurate and conservative PR evaluation.

Building upon the definitions of local robustness for individual input-label pairs (Def. 1, 2, 3), we extend the concept to a global robustness metric that quantifies the overall robustness of a classifier across the entire data distribution.

Definition 4 (Global Robustness) Consider input-label pairs $(\mathbf{x}, y) \sim D$, and let $\mathfrak{S}(\mathbf{x}, y)$ denote the point-wise robustness metric. The global robustness over the distribution D is defined as

$$\mathcal{G}(D) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim D} [\mathfrak{S}(\mathbf{x}, y)]. \quad (4)$$

For simplicity and to avoid ambiguity, we omit the distribution D from the notation \mathcal{G} and let \mathcal{G}_{AR} , \mathcal{G}_{PR} , and $\mathcal{G}_{\text{NPPR}}$ denote the respective global metrics for AR, PR, and NPPR.

Proposition 1 Considering AR, PR, and NPPR as defined in Def. 1, 2, and 3, and binary loss function for AR, let \mathcal{G}_{AR} , \mathcal{G}_{PR} , and $\mathcal{G}_{\text{NPPR}}$ denote their corresponding global robustness metrics. Given a perturbation distribution $\omega \in P_\varepsilon$ (either conditional or unconditional) for the perturbation ε , we have

$$\mathcal{G}_{\text{AR}} \leq \mathcal{G}_{\text{NPPR}} \leq \mathcal{G}_{\text{PR}}. \quad (5)$$

If we allow P_ε to be unrestricted, representing any family of distributions (including the Dirac delta measure), then the equality holds,

$$\mathcal{G}_{\text{AR}} = \mathcal{G}_{\text{NPPR}}. \quad (6)$$

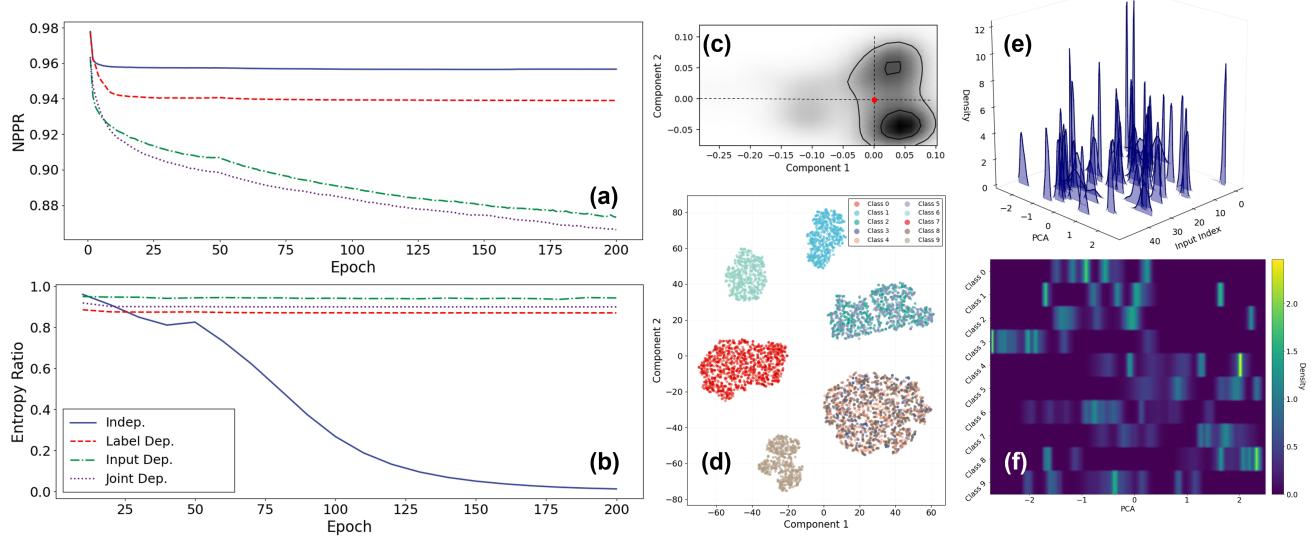


Figure 5. Training results of the proposed framework. Panels (a) and (b) show the training dynamics of the NPPR and the entropy ratio of the mixture proportions π_k over 200 epochs. We consider L_∞ -norm perturbations with radius $16/255$ under a Gaussian mixture model (GMM) with $K = 7$ modes. The solid blue curve corresponds to the independent case, where the same GMM parameters ϕ are used for all (\mathbf{x}, y) . The red dashed curve represents the label-dependent setting, while the green dash-dotted curve denotes the input-dependent case. The purple dotted curve illustrates the joint-dependency setting, in which the perturbation distribution depends on both \mathbf{x} and y . Panel (c) shows the PCA-projected contour of the perturbation distribution for the independent case (blue curve in Panels (a–b)). Panel (d) visualizes the t-SNE embeddings of perturbations for the label-dependent case (red dashed curve). Panel (e) presents PCA-based density plots for 50 randomly selected inputs under the input-dependent setting (green dash-dotted curve). Panel (f) displays a class-wise heatmap of perturbation densities for the jointly dependent case (purple dotted curve).

If P_ϵ includes only continuous distributions and the set of all adversarial perturbations is set of measure zero within \mathcal{B} , then the following strict inequality holds,

$$\mathcal{G}_{\text{AR}} < \mathcal{G}_{\text{NPPR}}. \quad (7)$$

Proposition 2 Reuse the condition in Prop. 1, and let \mathcal{G}^c and \mathcal{G}^u denote global robustness on conditional and unconditional perturbation distributions for AR and NPPR, respectively. Then we have

$$\mathcal{G}^c \leq \mathcal{G}^u. \quad (8)$$

Prop. 1 shows that the global NPPR metric serves as a more conservative measure compared to PR. Under extreme conditions, NPPR can be as low as AR. Prop. 2 compares the global robustness metrics for the conditional and unconditional cases and demonstrates that the conditional case yields lower robustness than the unconditional case. The detailed proofs can be found in Appendix 7.

4. NPPR Estimation via GMM

We formulate NPPR as an optimization problem parameterized via a GMM. The overall training pipeline of our

method is illustrated in Fig. 3. The pipeline begins by extracting intermediate features of the clean image \mathbf{x} from the target classifier, which are then fed into the perturbation generation process consisting of three components: (i) the *MLP heads*, (ii) the *GMM*, and (iii) a *bicubic up-sampling module*. Given the generated perturbations, a C&W-style loss (highlighted by the dashed block in Fig. 3) is computed using the logits produced by the classifier. We begin by introducing our objective function, followed by a detailed description of each component of our pipeline.

Objective Function Let \mathcal{U} denote the up-sampling module, which adjusts perturbations to match the input resolution while ensuring that their support lies within the prescribed perturbation budget. Consider a classifier h and a differentiable surrogate loss φ that relaxes the hard indicator loss 1. Given an input-label pair $(\mathbf{x}, y) \sim D$ and perturbations sampled from a GMM, our objective function is

$$\mathcal{L}(\phi) = \mathbb{E}_{(\mathbf{x}, y) \sim D} \left[\mathbb{E}_{\epsilon \sim \text{GMM}_\phi} \left[\varphi(h(\mathbf{x} + \mathcal{U}(\epsilon)), y) \right] \right], \quad (9)$$

Table 1. **Ablation settings of ResNet18 on CIFAR-10 under various perturbation dependency settings (%)**. $\widehat{\mathcal{G}}_{\text{NPPR}}$ denotes the estimator of the global NPPR metric (defined in Def. 3 and 4), while ER represents the entropy ratio of the mixture weights, indicating whether the distribution is dominated by a single mode. A smaller ER implies that the GMM is primarily governed by one dominant component. The lowest value of $\widehat{\mathcal{G}}_{\text{NPPR}}$ is highlighted in bold.

Config	Indep.		Label dep.		Input Dep.		Joint Dep.	
	$\widehat{\mathcal{G}}_{\text{NPPR}}$	ER(π)						
Base (K=3)	95.27	0.21	95.10	99.06	90.31	98.73	90.04	99.02
+ Learnable	95.85	91.35	93.64	99.21	90.26	94.98	89.72	98.94
+ #mode (K=7)	95.15	94.37	92.95	87.34	90.31	94.21	89.72	89.92
+ #mode (K=12)	95.53	92.43	92.68	76.59	90.02	79.54	89.50	73.83

Table 2. **Comparison of AR, PR, and NPPR for ResNet18 on CIFAR-10 under different perturbation budgets (%)**. We consider the L_∞ perturbation budgets with radii of 4/255, 8/255, and 16/255. For NPPR, different dependency settings are evaluated. PR is computed using a uniform perturbation distribution, and AR is obtained via PGD-20.

Estimator	Dependency	4/255	8/255	16/255
$\widehat{\mathcal{G}}_{\text{NPPR}}$	Indep.	99.26	98.21	95.15
	Label dep.	98.55	96.57	92.95
	Input dep.	98.56	96.53	90.31
	Joint dep.	97.91	95.56	89.71
$\widehat{\mathcal{G}}_{\text{PR}_{\text{Uniform}}}$		99.89	99.77	99.49
$\widehat{\mathcal{G}}_{\text{AR}_{\text{PGD}}}$		36.66	23.26	9.27

where ϕ is learnable parameters of GMM. In practice, since D is unknown, we use its empirical estimates

$$\widehat{\mathcal{L}}(\phi) = \frac{1}{MN} \sum_{i=1}^N \sum_{j=1}^M \left[\varphi(h(\mathbf{x}_i + \mathcal{U}(\boldsymbol{\varepsilon}_{i,j})), y_i) \right] \quad (10)$$

where $\boldsymbol{\varepsilon}_{i,j}$ denotes the j -th perturbation generated for image i by the GMM. We define the loss function φ as the logit margin between the ground-truth class and the most confident non-ground-truth class, with a hyperparameter κ controlling the margin scale. Let $\mathbf{x}' = \mathbf{x} + \mathcal{U}(\boldsymbol{\varepsilon})$ denote the AE. The loss function is

$$\varphi(h(\mathbf{x}'), y) = \text{softplus} \left(h_y(\mathbf{x}') - \max_{j \neq y} h_j(\mathbf{x}') + \kappa \right), \quad (11)$$

where $\text{softplus}(x) = \log(1 + e^x)$ lower-bounds the loss at zero, improves its smoothness near the origin, and behaves approximately linearly when the logit gap is large. We experimented with several variants of the loss function, and found that the *softplus* formulation yields the most stable

optimization trajectory. The loss function is highlighted by the dashed block in Fig. 3.

Our objective function can be interpreted as a global NPPR estimator that relaxes the hard 0–1 loss and models the perturbation distribution via a GMM. To capture different dependency structures between the perturbations and inputs, we reuse intermediary features from the classifier and pass them through an MLP to parameterize the GMM.

MLP heads We employ a two-layer MLP to model the dependency between the input and the perturbation distribution, using features extracted from the target classifier. As illustrated in Fig. 3, the features first pass through a shared fully connected layer, followed by batch normalization and a ReLU activation function. Subsequently, three independent fully connected layers are employed to produce the parameters of the GMM, including the mixture weights π_k , and the mean and covariance matrices μ_k and Σ_k for each component $k \in [K]$. The details of the GMM formulation are presented in the following subsection. In addition to modeling the dependency between the input \mathbf{x} and the perturbations, we also consider three other types of dependency structures illustrated in Fig. 4.

The first type is **(i) Independence**, where the perturbation distribution is entirely independent of the input–label pairs. This corresponds to scenarios such as additive Gaussian noise, external disturbances, and temperature drift, as discussed in Remark 1. **(ii) Label dependence:** As shown in Fig. 4 (a), the perturbation distribution depends on the ground-truth labels through a learnable label embeddings on mixture proportions. In this case, the noise is characterized by label-specific variations rather than input-dependent ones. **(iii) Input dependence:** The perturbation distribution depends solely on the input features extracted from the classifier. **(iv) Joint dependence:** As illustrated in Fig. 4 (b), the mixture weights π_k are linked to ground truth labels, while the means and covariances of the mixture components, μ_k and Σ_k , are conditioned on the input features.

GMM The GMM was first introduced by statistical researchers such as Titterington *et al.* [27] and later formalized by McLachlan and Peel [19] to model complex probability distributions. It has been theoretically shown that a GMM can approximate any continuous distribution arbitrarily well, given a sufficient number of components [2]. At a high level, a GMM models a distribution as a finite mixture of Gaussian components, each capturing a local mode of the underlying data distribution. The probabilistic graphical model of the GMM is illustrated in Fig. 3, where the mixture consists of K components (modes). The generative process for the i -th sample is formally defined as

$$z_i \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K), \quad (12)$$

$$\varepsilon_i | z_i \sim \mathcal{N}(\boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}), \quad (13)$$

where we first sample the categorical random variable z_i , parameterized by the mixture proportions, and then draw a sample from the z_i -th Gaussian component. We adopt the GMM to model the perturbation distribution ω due to its simplicity and strong theoretical foundations.

Training a GMM typically relies on the Expectation–Maximization (EM) algorithm. However, in our case, there are no pre-defined perturbation samples available, making the EM procedure inapplicable. Moreover, since the sampling process in the GMM involves a discrete latent variable z_i , standard gradient-based optimization cannot be directly applied. To address this issue, we employ the Gumbel–Softmax trick [15, 18], which provides a differentiable approximation to categorical sampling. At a high level, the Gumbel–Softmax reparameterization smooths the discrete mixture weights, thereby enabling gradient-based end-to-end optimization. The details on Gumbel–Softmax trick can be found in Appendix 9.1.

Bicubic up-sampling Since input images typically reside in high-dimensional spaces, the covariance matrix of the perturbation distribution scales as $\mathcal{O}(d^2)$ with respect to the input dimension d , making direct computation infeasible. To address this, we perform perturbation modeling in a lower-dimensional feature space and map the resulting perturbations back to the input space using bicubic interpolation, which is computationally efficient and preserves spatial smoothness [10, 17]. Detailed interpolation formulas are provided in the Appendix 8.

To ensure that the support of the perturbation distribution lies within the prescribed L_p -norm ball, we apply a scaled tanh mapping multiplied by the perturbation budget γ , a common constraint mechanism in the AR literature [6, 14].

5. Experiments

Experimental setup We conduct experiments on CIFAR-10, CIFAR-100, and Tiny ImageNet using ResNet18/50,

WideResNet50, and VGG16 as the base classifiers. Our GMM-based perturbation generator is trained with different numbers of mixture components ($K \in \{3, 7, 12\}$) under four conditioning strategies: independent, label-dependent, input-dependent, and jointly dependent. We evaluate performance under L_∞ -norm perturbation budget with radii $\varepsilon \in \{4/255, , 8/255, , 16/255\}$. The model is optimized using the C&W margin loss for 50 epochs with Adam ($\text{lr} = 5 \times 10^{-4}$). All experiments are implemented in Python 3.10 and PyTorch 2.5.1 on two NVIDIA RTX 3090 GPUs. Additional experimental details are provided in Appendix 9.

5.1. Learning Dynamics and Distribution Visualization

Fig. 5 illustrates the training dynamics (panels (a–b)) and the resulting distributions (panels (c–f)) of the proposed pipeline. For each training epoch, we record the running $\widehat{\mathcal{G}}_{\text{NPPR}}$ and the entropy ratio (ER) which measures the extent of mode dominance in a GMM, with lower values reflecting a collapse toward a single dominant component (see detail in Appendix 9.2). As shown in panel (a), the training process gradually yields a perturbation distribution under which the computed PR decreases over time. Moreover, as the dependency structure becomes more complex, the NPPR value further decreases (the jointly dependent case exhibits the lowest value, while the independent case remains the highest). This observation supports our theoretical analysis in Prop. 2.

As shown in panel (b), the independent case exhibits a clear collapse to a single dominant mode, indicating that the optimization tends to converge to a local minimum rather than exploring multiple mixture components.

Panel (c) presents the perturbation distribution for the independent case, projected onto the first two principal components using PCA. The results show that the learned distribution allocates a greater portion of its probability mass away from the center (marked by the red dot), indicating increased diversity that encourages exploration of decision boundaries.

Panel (d) visualizes the label-dependent perturbation distribution using t-SNE. The result reveals an interesting structural pattern. Although the distribution is conditioned on the ground-truth labels (10 in total), only 6 distinct clusters are formed. Some classes, such as class 7, are well separated, whereas others, e.g., classes 3 and 4, exhibit strong overlap, indicating shared perturbation characteristics. Fig. 6 in Appendix 9 presents the t-SNE visualization for the jointly dependent case, where samples from all classes are clearly disentangled.

Panel (e) randomly samples 50 CIFAR-10 inputs and applies PCA to their perturbations to visualize the dominant variation directions. Although an individual input may still

Table 3. **Performance across datasets and models (%)**. Mean over 10 runs (std in parentheses). Results are obtained under joint dependency with 7 mixture modes and an L_∞ perturbation radius of 16/255.

Dataset	Model	$\hat{\mathcal{G}}_{\text{NPPR}}$	$\hat{\mathcal{G}}_{\text{PR}_{\text{Gaussian}}}$	$\hat{\mathcal{G}}_{\text{PR}_{\text{Uniform}}}$	$\hat{\mathcal{G}}_{\text{AR}_{\text{PGD}}}$	$\hat{\mathcal{G}}_{\text{AR}_{\text{CW}}}$	Acc.
CIFAR10	ResNet18	88.32(0.63)	99.47(0.15)	99.60(0.22)	9.27(0.96)	9.67(0.99)	86.42
	ResNet50	92.99(0.44)	99.18(0.20)	99.40(0.26)	8.57(0.77)	8.87(0.61)	90.90
	WRN50	91.68(0.42)	99.28(0.12)	99.49(30.91)	8.89(2.94)	8.95(2.96)	91.11
	VGG16	97.83(0.34)	99.47(0.28)	99.62(0.17)	8.79(0.85)	9.31(0.77)	92.22
CIFAR100	ResNet18	74.17(0.48)	97.99(0.53)	98.60(0.61)	3.85(0.76)	4.21(0.61)	61.00
	ResNet50	85.16(0.52)	96.93(0.47)	97.80(0.35)	3.86(0.61)	4.32(0.45)	70.18
	WRN50	78.26(0.73)	97.48(0.42)	98.30(0.45)	4.53(0.66)	4.77(0.58)	70.42
	VGG16	91.81(0.64)	95.57(0.56)	95.81(0.62)	4.43(0.56)	4.19(0.60)	71.52
TinyImageNet	ResNet18	60.16(0.83)	98.87(0.29)	99.17(0.36)	2.76(0.50)	2.71(0.63)	57.65
	ResNet50	74.35(0.78)	98.94(0.19)	99.21(0.26)	3.56(0.66)	3.62(0.53)	69.42
	WRN50	64.66(0.77)	98.42(0.24)	98.79(0.32)	2.32(0.35)	2.60(0.39)	72.80
	VGG16	78.08(0.75)	98.96(0.28)	99.23(0.33)	4.69(0.75)	4.24(0.56)	64.14

exhibit mode collapse, which is consistent with the mixture-weight bar plots in Fig. 8 (Appendix 10), different inputs activate distinct perturbation patterns. This diversity in perturbations results in substantially lower NPPR values.

Panel (f) shows the heatmap for the jointly dependent case, where each band represents the perturbation distribution over 100 randomly selected inputs within a given class. Compared to the label-dependent case, the jointly dependent formulation offers greater flexibility in the distribution of each class (cf. Fig. 9 in Appendix 10).

5.2. Ablation Study

Tab. 1 presents the ablation study of our training pipeline on CIFAR-10 using ResNet18. As shown, except for the independent case, all models achieve improved performance, reflected by lower estimated $\hat{\mathcal{G}}_{\text{NPPR}}$, when a learnable up-sampler is used. The performance generally increases as the number of mixture components grows. The independent case behaves differently, which is due to its training instability.

5.3. Evaluation Across Models, Datasets, and Radii

Tab. 2 presents the results of NPPR, PR, and AR under different perturbation radii. As shown, across all cases, smaller radii correspond to higher evaluation result, indicating fewer AEs within the constrained perturbation region.

Tab. 3 summarizes and compares the global NPPR, PR (under Gaussian and uniform perturbation distributions), AR (evaluated with PGD-20 and C&W attacks), and clean accuracy across different DL models and datasets. As shown, all models exhibit the lowest AR performance under both PGD-20 and C&W attacks, while achieve the highest PR performance under both Gaussian and uniform perturbation distributions. Their NPPR performance consistently

lies between AR and general PR, confirming Prop. 1 that NPPR provides a more conservative and reliable robustness estimate. Our estimator further shows that NPPR performance deteriorates for the same model architecture as the underlying classification task becomes more challenging, e.g., $\hat{\mathcal{G}}_{\text{NPPR}}$ for ResNet18 drops from 88.32 on CIFAR-10 to 74.17 on CIFAR-100 and 60.61 on TinyImageNet.

6. Conclusion

To bridge the gap caused by the unrealistic reliance on predefined perturbation distributions in existing PR evaluations, we formally introduce Non-Parametric Probabilistic Robustness (NPPR). By directly learning from data, we further design a training pipeline and a GMM-based estimator for NPPR, consisting of: (i) multilayer perceptron (MLP) heads for modeling dependency structures, (ii) a Gaussian Mixture Model (GMM) for representing the perturbation distribution, and (iii) a bicubic up-sampling module for mapping perturbations to the original input resolution. We analyze the effects of different dependency structures and observe that stronger forms of dependency consistently yield lower NPPR values. Experiments across multiple datasets and model architectures demonstrate that our NPPR approach effectively learns a conservative perturbation distribution, offering a more practical solution given that the true perturbation distribution underlying PR can never be known in reality. Furthermore, we provide a theoretical analysis that clarifies the relationships among AR, PR, and NPPR under both input-dependent and input-independent perturbations. In summary, we believe NPPR represents an important stepping stone toward making PR more implementable in real-world applications.

References

- [1] Teodora Baluta, Zheng Leong Chua, Kuldeep S Meel, and Prateek Saxena. Scalable quantitative verification for deep neural networks. In *ICSE'21*, pages 312–323, 2021. [1](#) [2](#)
- [2] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006. [7](#)
- [3] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE symposium on security and privacy (sp)*, pages 39–57, 2017. [3](#)
- [4] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdip Mukhopadhyay. A survey on adversarial attacks and defences. *CAAI Transactions on Intelligence Technology*, 6(1):25–45, 2021. [1](#)
- [5] Ashutosh Chaudhary, Nikhil Agrawal, Kavya Barnwal, Keerat K Guliani, and Pramod Mehta. Universal adversarial perturbations: A survey. *arXiv preprint arXiv:2005.08087*, 2020. [1](#)
- [6] Jiefeng Chen, Xi Wu, Yang Guo, Yingyu Liang, and Somesh Jha. Towards evaluating the robustness of neural networks learned by transduction. In *International Conference on Learning Representations*, 2022. [7](#) [3](#)
- [7] Nicolas Couellan. Probabilistic robustness estimates for feed-forward neural networks. *Neural networks*, 142:138–147, 2021. [1](#) [2](#)
- [8] Elvis Dohmatob and Alberto Bietti. On the (non-) robustness of two-layer neural networks in different learning regimes. *arXiv preprint arXiv:2203.11864*, 2022. [2](#)
- [9] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. [2](#)
- [10] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *European conference on computer vision*, pages 391–407. Springer, 2016. [7](#) [3](#)
- [11] Ian J Goodfellow. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations ICLR 2015*. International Conference on Learning Representations, 2015. [1](#)
- [12] Wei Huang, Xingyu Zhao, Gaojie Jin, and Xiaowei Huang. Safari: Versatile and efficient evaluations for robustness of interpretability. In *IEEE/CVF International Conference on Computer Vision (ICCV'23)*, pages 1988–1998, 2023. [4](#)
- [13] Xiaowei Huang, Daniel Kroening, Wenjie Ruan, and et al. A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence, and interpretability. *Computer Science Review*, 37:100270, 2020. [1](#) [2](#)
- [14] Zhichao Huang and Tong Zhang. Black-box adversarial attack with transferable model-based embedding. In *International Conference on Learning Representations*, 2020. [7](#) [3](#)
- [15] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. [7](#) [3](#)
- [16] Herman Kahn and Theodore E Harris. Estimation of particle transmission by random sampling. *National Bureau of Standards applied mathematics series*, 12:27–30, 1951. [4](#)
- [17] Robert Keys. Cubic convolution interpolation for digital image processing. *IEEE transactions on acoustics, speech, and signal processing*, 29(6):1153–1160, 2003. [7](#) [3](#)
- [18] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016. [7](#) [3](#)
- [19] Geoffrey J McLachlan and David Peel. *Finite mixture models*. John Wiley & Sons, 2000. [7](#)
- [20] Mark Huasong Meng, Guangdong Bai, Sin Gee Teo, Zhe Hou, Yan Xiao, Yun Lin, and Jin Song Dong. Adversarial robustness of deep neural networks: A survey from a formal verification perspective. *IEEE Transactions on Dependable and Secure Computing*, 2022. [1](#)
- [21] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *IEEE European symposium on security and privacy (EuroS&P)*, pages 372–387. IEEE, 2016. [1](#)
- [22] Mikhail Pautov, Nurislam Tursynbek, Marina Munkhoeva, Nikita Muravev, Aleksandr Petushko, and Ivan Oseledets. Cc-cert: A probabilistic approach to certify general robustness of neural networks. In *AAAI'22*, pages 7975–7983, 2022. [1](#) [2](#)
- [23] Alexander Robey, Luiz Chamon, George J Pappas, and Hamed Hassani. Probabilistically robust learning: Balancing average and worst-case performance. In *International Conference on Machine Learning*, pages 18667–18686. PMLR, 2022. [1](#) [2](#) [4](#)
- [24] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *In Proc. of 2nd Int. Conf. on Learning Representations*, 2014. [1](#)
- [25] Karim Tit, Teddy Furon, and Mathias Rousset. Efficient statistical assessment of neural network corruption robustness. *NeurIPS'21*, 34:9253–9263, 2021. [1](#) [2](#)
- [26] Karim Tit, Teddy Furon, and Mathias Rousset. Gradient-informed neural network statistical robustness estimation. In *Proc. of The 26th Int. Conf. on Artificial Intelligence and Statistics*, pages 323–334. PMLR, 2023. [1](#) [2](#) [4](#)
- [27] David Michael Titterington, Adrian FM Smith, and Udi E Makov. Statistical analysis of finite mixture distributions. (*No Title*), 1985. [7](#)
- [28] Benjie Wang, Stefan Webb, and Tom Rainforth. Statistically robust neural network classification. In *Uncertainty in Artificial Intelligence*, pages 1735–1745. PMLR, 2021. [1](#) [2](#) [4](#)
- [29] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1905–1914, 2021. [2](#)
- [30] Zheng Wang, Geyong Min, and Wenjie Ruan. The implicit bias of gradient descent toward collaboration between layers: A dynamic analysis of multilayer perceptions. *Advances in Neural Information Processing Systems*, 37:74868–74898, 2024. [2](#)

- [31] Stefan Webb, Tom Rainforth, Yee Whye Teh, and M Pawan Kumar. A statistical approach to assessing neural network robustness. In *Int. Conf. on Learning Representations*, 2019. [1](#), [2](#), [3](#), [4](#)
- [32] Lily Weng, Pin-Yu Chen, Lam Nguyen, Mark Squillante, Akhilan Boopathy, Ivan Oseledets, and Luca Daniel. Proven: Verifying robustness of neural networks with a probabilistic approach. In *Int. Conf. on Machine Learning*, pages 6727–6736. PMLR, 2019. [2](#)
- [33] Tianle Zhang, Wenjie Ruan, and Jonathan E. Fieldsend. Proa: A probabilistic robustness assessment against functional perturbations. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2022, Grenoble, France, September 19–23, 2022, Proceedings, Part III*, page 154–170. Springer-Verlag, 2023. [1](#), [2](#), [3](#), [4](#)
- [34] Tianle Zhang, Yanghao Zhang, Ronghui Mu, Jiaxu Liu, Jonathan Fieldsend, and Wenjie Ruan. Prass: probabilistic risk-averse robust learning with stochastic search. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 559–567, 2024. [2](#), [4](#)
- [35] Yi Zhang, Yun Tang, Wenjie Ruan, Xiaowei Huang, Siddartha Khastgir, Paul Jennings, and Xingyu Zhao. Protip: Probabilistic robustness verification on text-to-image diffusion models against stochastic perturbation. In *European Conference on Computer Vision*, 2024. [1](#), [2](#), [3](#), [4](#)
- [36] Yi Zhang, Yuhang Chen, Zhen Chen, Wenjie Ruan, Xiaowei Huang, Siddartha Khastgir, and Xingyu Zhao. Adversarial training for probabilistic robustness. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1675–1685, 2025. [1](#)
- [37] Yi Zhang, Yuhang Chen, Zhen Chen, Wenjie Ruan, Xiaowei Huang, Siddartha Khastgir, and Xingyu Zhao. Adversarial training for probabilistic robustness. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1675–1685, 2025. [2](#), [4](#)
- [38] Yi Zhang, Zheng Wang, Chen Zhen, Wenjie Ruan, Qing Guo, Siddartha Khastgir, Carsten Maple, and Xingyu Zhao. Probabilistic robustness for free? revisiting training via a benchmark, 2025. [1](#), [2](#), [4](#)
- [39] Xingyu Zhao. Probabilistic robustness in deep learning: A concise yet comprehensive guide. *Adversarial Example Detection and Mitigation Using Machine Learning*, pages 1–13, 2025. [1](#), [3](#), [4](#)
- [40] Monty-Maximilian Zühlke and Daniel Kudenko. Adversarial robustness of neural networks from the perspective of lipschitz calculus: A survey. *ACM Computing Surveys*, 2024. [1](#), [2](#)

Non-Parametric Probabilistic Robustness: A Conservative Metric with Optimized Perturbation Distributions

Supplementary Material

7. Omitted Proofs

To facilitate readability, we restate Propositions 1 and 2, followed by their proofs in order.

Proposition 3 Considering AR, PR, and NPPR as defined in Def. 1, 2, and 3, and binary loss function for AR, let \mathcal{G}_{AR} , \mathcal{G}_{PR} , and $\mathcal{G}_{\text{NPPR}}$ denote their corresponding global robustness metrics. Given a perturbation distribution $\omega \in P_\epsilon$ (either conditional or unconditional) for the perturbation ϵ , we have

$$\mathcal{G}_{\text{AR}} \leq \mathcal{G}_{\text{NPPR}} \leq \mathcal{G}_{\text{PR}}. \quad (14)$$

If we allow P_ϵ to be unrestricted, representing any family of distributions (including the Dirac delta measure), then the equality holds,

$$\mathcal{G}_{\text{AR}} = \mathcal{G}_{\text{NPPR}}. \quad (15)$$

If P_ϵ includes only continuous distributions and the set of all adversarial perturbations is set of measure zero within \mathcal{B} , then the following strict inequality holds,

$$\mathcal{G}_{\text{AR}} < \mathcal{G}_{\text{NPPR}}. \quad (16)$$

Proof 1 Here we provide the proof of Proposition 1. The inequality $\mathcal{G}_{\text{NPPR}} \leq \mathcal{G}_{\text{PR}}$ follows directly from Definition 3. Thus, it remains to establish the inequality between AR and NPPR, namely $\mathcal{G}_{\text{AR}} \leq \mathcal{G}_{\text{NPPR}}$. We prove this inequality by showing that $\mathcal{G}_{\text{AR}} = \mathcal{G}_{\text{NPPR}}$ when no restrictions are imposed on P_ϵ , and for some specific restrictions imposed on P_ϵ , there exists $\mathcal{G}_{\text{AR}} < \mathcal{G}_{\text{NPPR}}$.

We begin with the equality. To this end, we first establish that $\mathcal{G}_{\text{AR}} \geq \mathcal{G}_{\text{NPPR}}$, and then show the reverse inequality $\mathcal{G}_{\text{AR}} \leq \mathcal{G}_{\text{NPPR}}$. Considering the conditional case and a binary loss function, let

$$\epsilon^* \in \arg \sup_{\epsilon \in \mathcal{B}} \mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}] \quad (17)$$

be the adversarial perturbation, and

$$\omega^*(\cdot | \mathbf{x}, y) = \arg \inf_{\omega \in P_\epsilon} \mathbb{E}_{\epsilon \sim \omega(\cdot | \mathbf{x}, y)} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}] \quad (18)$$

be optimal perturbation distribution of NPPR. Now, considering Dirac delta measure $\delta_{\epsilon^*} \in P_\epsilon$, we have

$$\mathcal{G}_{\text{AR}} = \mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y}] \quad (19)$$

$$= \mathbb{E}_D [\mathbb{E}_{\delta_{\epsilon^*}} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (20)$$

$$\geq \mathbb{E}_D [\inf_{\omega} \mathbb{E}_{\omega} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (21)$$

$$= \mathcal{G}_{\text{NPPR}} \quad (22)$$

Inversely, we have

$$\mathcal{G}_{\text{NPPR}} = \mathbb{E}_D [\mathbb{E}_{\omega^*} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]], \quad (23)$$

If $\mathbb{E}_{\omega^*} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}] < 1$, there must exist at least one $\epsilon^* \in \mathcal{B}$, such that $h(\mathbf{x} + \epsilon^*) \neq y$, therefore $\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y} = 0$, hence

$$\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y}] \leq \mathbb{E}_D [\mathbb{E}_{\omega^*} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]]. \quad (24)$$

Therefore, if we do not restrict P_ϵ , we have $\mathcal{G}_{\text{NPPR}} = \mathcal{G}_{\text{AR}}$. In case of P_ϵ represents continuous distributions, and there only exists distinct adversarial examples ϵ^* such that $\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y} = 0$ with probability of zero, then $\forall \omega \in P_\epsilon$ we have $\mathbb{E}_{\omega} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}] = 1$. Hence, we have

$$\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y}] < \mathbb{E}_D [\mathbb{E}_{\omega^*} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]]. \quad (25)$$

In the unconditional case, instead of the usual input-dependent PGD, we consider Universal Adversarial Attacks (UAEs) [5]. The following subsection will show that the usual input-dependent attack yields a lower value than that of the UAE. Let

$$\omega^* = \arg \inf_{\omega \in P_\epsilon} \mathbb{E}_D [\mathbb{E}_{\epsilon \sim \omega} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (26)$$

$$= \arg \inf_{\omega \in P_\epsilon} \mathbb{E}_{\epsilon \sim \omega} [\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (27)$$

The two expectations is exchangeable since ω is independent of input-label pairs. Let

$$\epsilon^* \in \arg \sup_{\epsilon \in \mathcal{B}} \mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon) \neq y}] \quad (28)$$

be an UAE. Similarly, we have

$$\mathcal{G}_{\text{AR}} = \mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon^*)=y}] \quad (29)$$

$$= \mathbb{E}_{\delta_{\epsilon^*}} [\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (30)$$

$$\geq \inf_{\omega} \mathbb{E}_{\omega} [\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (31)$$

$$= \mathcal{G}_{\text{NPPR}} \quad (32)$$

Now, we show that $\mathcal{G}_{\text{AR}} \leq \mathcal{G}_{\text{NPPR}}$. We have

$$\mathcal{G}_{\text{NPPR}} = \mathbb{E}_D [\mathbb{E}_{\omega^*} [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]] \quad (33)$$

$$= \mathbb{E}_{\omega^*} [\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\epsilon)=y}]]. \quad (34)$$

where ω^* is the optimal distribution derived from NPPR, and for any distribution ω , we have

$$\mathbb{E}_\omega [\mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y}]] \geq \inf_{\boldsymbol{\varepsilon}} \mathbb{E}_D [\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y}]. \quad (35)$$

Hence we have $\mathcal{G}_{\text{NPPR}} \geq \mathcal{G}_{\text{AR}}$. The proof of inequality $\mathcal{G}_{\text{NPPR}} > \mathcal{G}_{\text{AR}}$ is the same as the conditional case.

Now, we prove the Prop. 2.

Proposition 4 Reuse the condition in Prop. 1, and let \mathcal{G}^c and \mathcal{G}^u denote global robustness on conditional and unconditional perturbation distributions for AR, and NPPR, respectively. Then we have

$$\mathcal{G}^c \leq \mathcal{G}^u. \quad (36)$$

Proof 2 We first prove the AR case and consider the UAE as the unconditional case of AR. We have

$$\mathcal{G}_{\text{AR}}^c = \mathbb{E}_D \left[\sup_{\boldsymbol{\varepsilon}} \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon}) \neq y} \right] \quad (37)$$

$$= \mathbb{E}_D \left[\inf_{\boldsymbol{\varepsilon}} \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y} \right]. \quad (38)$$

Since $\forall \boldsymbol{\varepsilon}_0 \in \mathcal{B}$ and $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$, $\inf_{\boldsymbol{\varepsilon}} \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y} \leq \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon}_0)=y}$, hence $\forall \boldsymbol{\varepsilon}_0 \in \mathcal{B}$

$$\mathbb{E}_D \left[\inf_{\boldsymbol{\varepsilon}} \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y} \right] \leq \mathbb{E}_D \left[\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon}_0)=y} \right]. \quad (39)$$

Therefore,

$$\mathbb{E}_D \left[\inf_{\boldsymbol{\varepsilon}} \mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y} \right] \leq \inf_{\boldsymbol{\varepsilon}} \mathbb{E}_D \left[\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y} \right] = \mathcal{G}_{\text{AR}}^u. \quad (40)$$

Following the same logic, we show that.

$$\mathcal{G}_{\text{NPPR}}^c = \mathbb{E}_D \left[\inf_{\omega} \mathbb{E}_{\omega} [\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y}] \right] \quad (41)$$

$$\leq \inf_{\omega} \mathbb{E}_D [\mathbb{E}_{\omega} [\mathbf{1}_{h(\mathbf{x}+\boldsymbol{\varepsilon})=y}]] \quad (42)$$

$$= \mathcal{G}_{\text{NPPR}}^u. \quad (43)$$

8. Details on Bicubic Up-sampling

Since input images typically reside in high-dimensional spaces, the covariance matrix of the perturbation distribution becomes prohibitively large, scaling as $\mathcal{O}(d^2)$ with respect to the input dimension d . This quadratic growth renders both storage and computation infeasible when the input dimension is large, as in modern image datasets. To mitigate this issue, we perform the perturbation modeling in a lower-dimensional space, reducing computational overhead, and subsequently map the perturbations back to the input space using *bicubic interpolation*. This approach is computationally efficient while preserving the spatial smoothness of the perturbations, which has been studied in the robustness-related literature [9, 29].

Table 4. Configuration summary for NPPR estimation.

Setting	Value
Model / Architecture	
Up-sampler	bicubic interpolation
norm	ℓ_∞
ϵ	{4/255, 8/255, 16/255}
GMM Parameters	
Initialization	uniform
Modes K	{3, 7, 12}
Latent dim.	{128, 256}
Covariance type	full
Hidden dim.	{256, 512}
Label emb. dim.	{64, 128}
Label emb. norm.	TRUE
Training Hyperparameters	
Epochs	50
Learning rate	{ 5×10^{-4} , 2×10^{-2} }
LR warmup epochs	20
LR min	2×10^{-6}
Loss type	C&W
κ	1
Samples per input	32
Annealing Schedule	
T_π (init \rightarrow final)	3.0 \rightarrow 1.0
T_μ (init \rightarrow final)	3.0 \rightarrow 1.0
T_σ (init \rightarrow final)	1.5 \rightarrow 1.0
T_{shared} (init \rightarrow final)	1.5 \rightarrow 1.0
Gumbel anneal	TRUE
Gumbel temp (init \rightarrow final)	1.0 \rightarrow 0.1

Our bicubic up-sampling is composed of a linear mapping and a bicubic interpolation module. For a given pixel position (x, y) in the upsampled image, bicubic interpolation estimates its intensity as a weighted sum of the 4×4 neighboring pixels in the original image:

$$I'(x, y) = \sum_{m=-1}^2 \sum_{n=-1}^2 w(m, x) w(n, y) I(i+m, j+n), \quad (44)$$

where $I(i+m, j+n)$ denotes the neighboring pixel values and $w(\cdot, \cdot)$ represents the interpolation weights determined by a cubic convolution kernel. The one-dimensional cubic kernel $w(a)$ is defined as a piecewise cubic polyno-

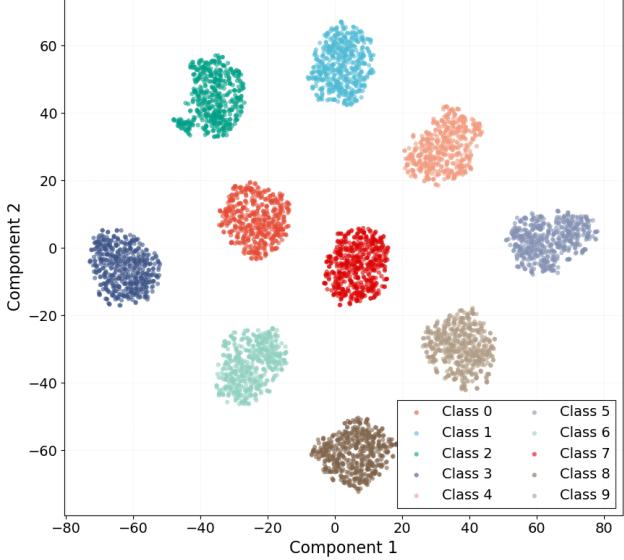


Figure 6. **The t-SNE plot for the jointly dependent case.** We additionally visualize the jointly conditioned model for ResNet18 on CIFAR-10. With the added dependence on inputs, the perturbation distributions for different classes become fully disentangled.

mial [17]:

$$w(a) = \begin{cases} (1.5)|a|^3 - 2.5|a|^2 + 1, & \text{if } |a| < 1, \\ -0.5|a|^3 + 2.5|a|^2 - 4|a| + 2, & \text{if } 1 \leq |a| < 2, \\ 0, & \text{otherwise.} \end{cases} \quad (45)$$

This kernel ensures smoothness and locality, producing continuous first derivatives while limiting interpolation to the 4×4 neighborhood around (i, j) . Bicubic up-sampling has been widely adopted as a baseline in image super-resolution and restoration tasks [10].

To ensure that the support of the perturbation distribution lies within the prescribed L_p -norm ball, we apply the mapping g_B , defined as

$$g_B = \gamma \tanh(\cdot), \quad (46)$$

which is a commonly used constraint mechanism in robustness literature [6, 14].

9. Detailed Experiment Settings

We provide the detailed experimental configurations in Tab. 4. For the independent-perturbation setting with a fixed up-sampler, we adopt a different training strategy from the other cases because this setting is substantially harder to optimize. Specifically, we use a larger learning rate of 2×10^{-2} with a cosine cyclical scheduler and a 20-epoch warm-up, which yields the most stable training trajectory.

For all other dependency settings, including those with a learnable up-sampler, we use a fixed learning rate of 5×10^{-4} . Except for the runs shown in Fig. 5, which are trained for 200 epochs for visualization purposes, all reported results are trained for 50 epochs.

Annealing Schedule We adopt an annealing strategy for the mixture weights as well as the parameters of each mixture-component distribution. This is motivated by the substantial imbalance in the number of parameters associated with the mixture weights, means, and covariance matrices. For example, when using $K = 3$, a latent dimension of 128, and a hidden dimension of 256, the mixture weights require only 3×256 parameters, whereas the mean and covariance heads require 128×256 and $128^2 \times 256$ parameters, respectively. Such a disparity can cause optimization to be dominated by the larger parameter groups, leading to suboptimal local minima. To mitigate this imbalance, we apply annealing to stabilize training and prevent premature convergence to poor solutions.

9.1. Gumbel softmax trick

Training a Gaussian Mixture Model (GMM) within a gradient-based framework requires differentiating through the discrete mixture-selection variable. Specifically, for each perturbation sample ϵ_i , a categorical latent variable $z_i \in \{1, \dots, K\}$ determines which Gaussian component generates the sample. Directly sampling $z_i \sim \text{Cat}(\pi_1, \dots, \pi_K)$ is non-differentiable, preventing back-propagation. To overcome this limitation, we adopt the Gumbel–Softmax (also known as the Concrete) relaxation [15, 18], which provides a differentiable approximation to categorical sampling.

The trick relies on the Gumbel perturbation property: if g_k are i.i.d. samples from $\text{Gumbel}(0, 1)$, then

$$z = \arg \max_k (\log \pi_k + g_k) \quad (47)$$

is exactly distributed as a categorical random variable with probabilities $\{\pi_k\}_{k=1}^K$. Instead of taking the non-differentiable $\arg \max$, Gumbel–Softmax introduces a temperature-controlled softmax relaxation:

$$\tilde{z}_k = \frac{\exp((\log \pi_k + g_k)/\tau)}{\sum_{j=1}^K \exp((\log \pi_j + g_j)/\tau)}, \quad k = 1, \dots, K, \quad (48)$$

where $\tau > 0$ is a temperature parameter. When $\tau \rightarrow 0$, the distribution becomes increasingly “one-hot,” recovering a true categorical sample; when τ is larger, the distribution is smoother, enabling stable gradients.

The reparameterized mixture selection is therefore given by the continuous vector

$$\tilde{z} = (\tilde{z}_1, \dots, \tilde{z}_K), \quad (49)$$

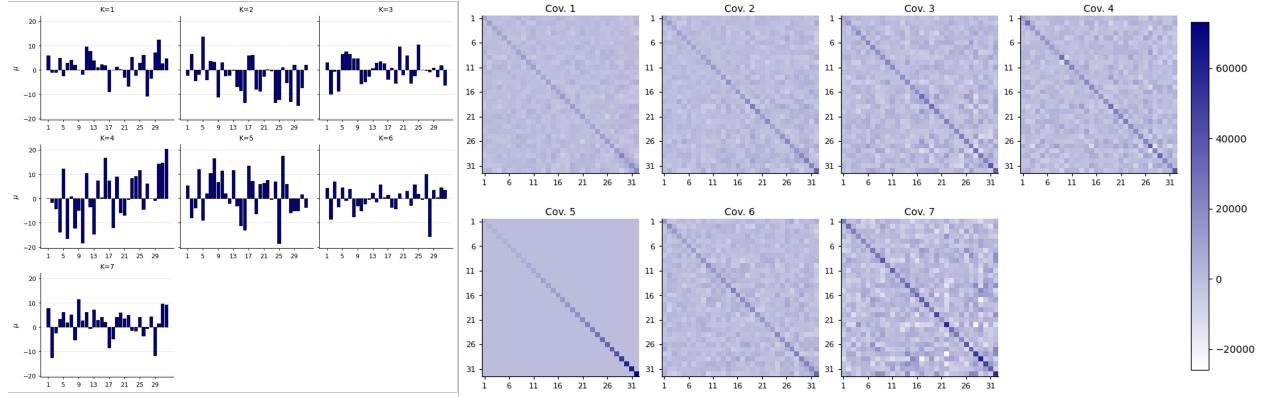


Figure 7. **Bar plot and heatmap of mixture component means and covariances.** For the input-dependent case on ResNet18 (CIFAR-10), we randomly select one input and visualize the GMM parameters after reducing the feature dimension to 32 using PCA. Specifically, we display a bar plot of the mixture means and a heatmap of the corresponding covariance matrices.

which lies in the probability simplex and is fully differentiable with respect to the mixture weights π_k . This relaxed one-hot vector replaces the discrete indicator and allows the GMM sample to be expressed as:

$$\varepsilon_i = \sum_{k=1}^K \tilde{z}_k \mu_k + \sum_{k=1}^K \tilde{z}_k \Sigma_k^{1/2} \xi_k, \quad (50)$$

where $\xi_k \sim \mathcal{N}(0, I)$ is an auxiliary noise variable. Because all operations are differentiable, the entire perturbation generation process is trainable via standard backpropagation.

During training, we anneal the temperature τ from a higher initial value to a smaller final value, which encourages exploration early on and progressively sharpens the mixture assignments. This annealing strategy stabilizes optimization and prevents premature distribution collapse.

9.2. Entropy Ratio

Entropy Ratio (ER) quantifies the degree of mode dominance in a Gaussian Mixture Model (GMM). It measures how evenly the mixture weights $\pi = (\pi_1, \dots, \pi_K)$ are distributed across the K components. Lower ER values indicate that the probability mass is concentrated on a single dominant mode (i.e., mode collapse), whereas values closer to 1 suggest a more uniform mixture distribution.

Formally, ER is defined as

$$\text{ER}(\pi) = \frac{H(\pi)}{\log K} = \frac{-\sum_{k=1}^K \pi_k \log \pi_k}{\log K}, \quad (51)$$

where $H(\pi)$ denotes the Shannon entropy of the mixture weights, and $\log K$ is the maximum possible entropy for a K -component mixture. This normalization ensures that $\text{ER} \in [0, 1]$, allowing consistent comparison across different values of K .

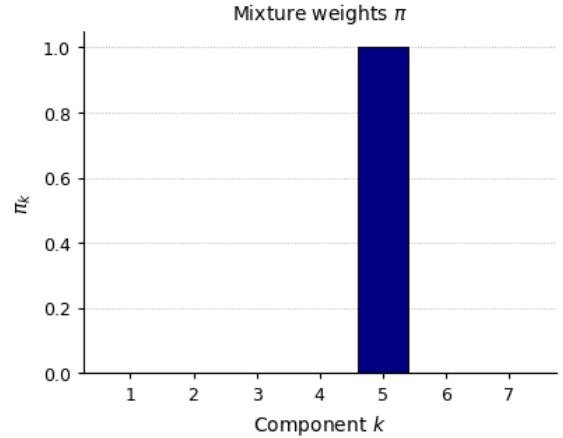


Figure 8. **Bar plot of mixture proportions.** We visualize the mixture proportions for the same input and model used in Fig. 7 by plotting the corresponding bar chart.

10. Additional Experiments

Here, we provide additional results that further illustrate the characteristics of the learned perturbation distributions.

Fig. 6 presents the t-SNE visualization corresponding to Fig. 5 panel (d), but using the joint-dependence structure instead of label dependence. As shown, perturbation samples from all classes become clearly disentangled, indicating that joint dependence yields a substantially more diverse and well-separated perturbation distribution than label dependence.

Fig. 7 and 8 present bar plots of the mixture means after applying PCA (reduced to 32 dimensions), heatmaps of the covariance matrices, and bar plots of the mixture proportions. The results show that only a single mixture component remains active, thereby dominating the distribution.

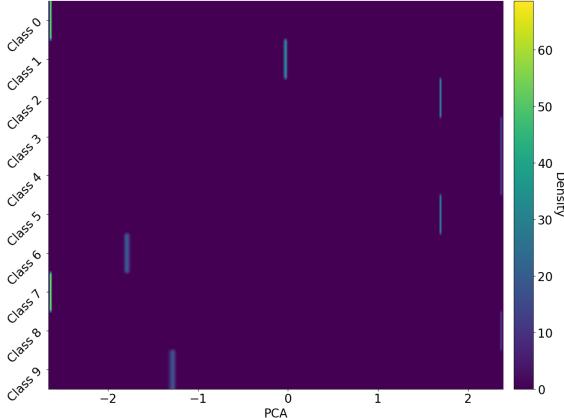


Figure 9. Class-wise heatmap of perturbation densities for the label-dependent case. The heatmap is generated from the learned distribution of ResNet on CIFAR-10, using the same experimental setup as panel (f) in Fig. 5.

The covariance heatmaps further reveal that this dominant component exhibits strong diagonal values with minimal off-diagonal structure, indicating limited correlation across dimensions.

Fig. 9 shows the heatmap of perturbation distributions grouped by labels. Compared with Fig. 5 panel (f), the joint-dependent case exhibits substantially greater diversity within each label group, indicating a more varied perturbation distribution.

Tab. 5 shows the extended experiments of our proposed pipeline using ResNet18 on CIFAR-10. As indicated, the results match our observations in the ablation study in Tab. 1. In addition, we include a setting where we directly optimize in the input space without the up-sampling module. Due to the large input dimensionality, this greatly increases the number of parameters and computational cost, and leads to a very unstable training trajectory. Without an up-sampler, it also hardly generalizes to high-resolution images, e.g., $3 \times 224 \times 224$ in ImageNet. Hence, we exclude it from our main experiments, though it has a better performance.

In the table, in addition to the estimated NPPR on the test set, $\hat{G}_{\text{NPPR}_{\text{test}}}$, and the entropy ratio, we also include the estimated NPPR on the training set and report the maximum, minimum, and standard deviation of the mixture proportions π . One interesting observation is that the estimated NPPR on the training set under the independent and label-dependent settings is higher than the value on the test data. Typically, this value should be lower on the training set and higher on the test set, with the difference reflecting the generalization gap. However, due to training instability in the independent case, the value on the training set becomes lower than expected. Moreover, the generalization gap between the train and test results for both the input-dependent

and joint-dependent settings becomes larger when the bicubic up-sampling module is removed. It shows the benefit of using the up-sampling module, namely a smaller generalization gap.

Table 5. **Extended Experimental Results of ResNet18 on CIFAR-10 (%)**. We also tested a variant that removes the up-sampler and optimizes perturbations directly in input space. Note that this setting is train unstable since it requires matching the full input dimensionality. It is also computational expensive, though it can yield a better performance. Hence, We exclude it from the main results.

Config.		$\widehat{\mathcal{G}}_{\text{NPPR}_{\text{test}}}$	$\widehat{\mathcal{G}}_{\text{NPPR}_{\text{train}}}$	ER(π)	Max(π)	Min(π)	Std.(π)
(I) Independent Perturbations							
3	None	94.25	95.07	0.000431	99.9960	0.0016	57.73
3	Non-trainable	95.26	95.76	0.002135	99.9766	0.0110	57.71
3	Trainable	95.85	96.30	0.913476	54.43	20.94	18.37
7	None	94.48	95.11	0.000290	99.9956	0.0006	37.79
7	Non-trainable	95.06	95.71	0.000935	99.9840	0.0013	37.79
7	Trainable	95.15	95.65	0.943673	31.49	8.38	7.93
12	None	94.59	95.43	0.000287	99.9946	0.0003	28.87
12	Non-trainable	95.06	95.74	0.000757	99.9844	0.0010	28.86
12	Trainable	95.53	95.88	0.924282	20.11	3.83	5.83
(II) Input-Dependent Perturbations (x-dependent)							
3	None	86.14	79.28	0.9216	52.25	18.87	17.13
3	Non-trainable	90.31	88.92	0.9873	40.95	27.59	6.87
3	Trainable	90.26	87.97	0.9498	46.63	19.91	13.36
7	None	85.60	77.33	0.8512	30.06	2.05	11.32
7	Non-trainable	90.20	89.22	0.8991	29.46	2.63	9.45
7	Trainable	90.31	87.33	0.9421	24.38	4.61	7.16
12	None	85.13	75.04	0.9404	17.55	1.45	4.62
12	Non-trainable	89.41	89.98	0.7829	29.30	0.00	8.90
12	Trainable	90.02	87.26	0.7954	20.22	0.00	7.95
(III) Label-Dependent Perturbations (y-dependent)							
3	None	90.48	92.23	0.9885	40.89	28.75	6.59
3	Non-trainable	95.10	95.99	0.9906	40.14	29.15	5.94
3	Trainable	93.64	94.54	0.9921	39.64	30.11	5.46
7	None	89.30	90.69	0.7359	40.08	0.0001	14.85
7	Non-trainable	94.03	94.92	0.7320	40.27	0.0117	15.05
7	Trainable	92.95	93.90	0.8734	30.27	0.1151	9.82
12	None	89.22	90.77	0.5458	50.23	0.000034	14.72
12	Non-trainable	94.50	95.44	0.7068	20.35	0.0113	9.34
12	Trainable	92.68	94.01	0.7659	20.22	0.0095	8.25
(IV) Joint-Dependent Perturbations (x, y-dependent)							
3	None	84.82	78.53	0.9369	50.43	20.51	15.41
3	Non-trainable	90.04	89.24	0.9902	40.31	29.30	6.07
3	Trainable	89.72	87.13	0.9894	40.60	29.24	6.31
7	None	84.10	72.09	0.9689	22.34	8.91	5.43
7	Non-trainable	89.22	87.83	0.8786	29.22	0.0001	9.34
7	Trainable	89.72	86.62	0.8992	20.40	0.0371	7.89
12	None	84.19	75.93	0.7610	19.96	0.0004	8.28
12	Non-trainable	89.42	89.25	0.7376	30.23	0.0008	9.41
12	Trainable	89.50	86.71	0.7383	30.60	0.0191	9.49