

# Robustness-via-Synthesis: Robust Training with Generative Adversarial Perturbations

Inci M. Baytaş, and Debayan Deb

**Abstract**—Upon the discovery of adversarial attacks, robust models have become obligatory for deep learning-based systems. Adversarial training with first-order attacks has been one of the most effective defenses against adversarial perturbations to this day. The majority of the adversarial training approaches focus on iteratively perturbing each pixel with the gradient of the loss function with respect to the input image. However, the adversarial training with gradient-based attacks lacks diversity and does not generalize well to natural images and various attacks. This study presents a robust training algorithm where the adversarial perturbations are automatically synthesized from a random vector using a generator network. The classifier is trained with cross-entropy loss regularized with the optimal transport distance between the representations of the natural and synthesized adversarial samples. Unlike prevailing generative defenses, the proposed one-step attack generation framework synthesizes diverse perturbations without utilizing gradient of the classifier’s loss. Experimental results show that the proposed approach attains comparable robustness with various gradient-based and generative robust training techniques on CIFAR10, CIFAR100, and SVHN datasets. In addition, compared to the baselines, the proposed robust training framework generalizes well to the natural samples. Code and trained models will be made publicly available.

**Index Terms**—Adversarial robustness, adversarial training, adversarial attacks synthesis, optimal transport.

## I. INTRODUCTION

DEEP neural networks, in particular, convolutional neural networks (CNNs) have become a cornerstone for representation learning in various computer vision applications. Remarkable state-of-the-art performances of CNN architectures have been reported for various benchmark datasets in literature [1]. On the other hand, Szegedy *et al.* [2] unraveled an unprecedented vulnerability of CNNs to specifically crafted, but imperceptible, perturbations known as adversarial attacks. Upon this discovery, various studies [3]–[7] showed that it is possible to generate adversarial perturbations of varying strengths using a target, or a surrogate deep model.

The adversarial attacks can be classified into three categories, such as white-box, gray-box, and black-box attacks [8]. The white-box attacks require the full knowledge of the target model [8]. Projected Gradient Descent (PGD) [9] is accepted as one of the most powerful white-box attacks. The PGD method iteratively computes the first-order gradient of the

target model with respect to the input. Whereas, the gray-box attacks have access to the architecture of the target model however the learned model parameters are not available [8]. On the other hand, the black-box methods are only allowed to make queries to the model [8]. The presence of various types of adversarial attacks indicates that the adversarial samples may emerge as a security risk for many deep learning-based systems, such as self-driving cars [10], face recognition systems [11], and healthcare [12]. Therefore, it is now imperative to safeguard deep networks against adversarial perturbations.

In a standard supervised training setting, CNNs are trained with a set of natural images which lack representation of potential adversarial samples. Thus, a naturally trained deep classifier is prone to misclassify the adversarial samples. To alleviate the vulnerability of the CNNs against adversarial samples, adversarial defense methods have been developed. The state-of-the-art defense techniques pose the adversarial robustness as a generalization problem and try to regularize the deep model training by augmenting the training set with adversarial samples. Such techniques are known as *adversarial training* [9], [13]. In the most common adversarial training methods [9], [13], the adversarial samples are generated by a single type of adversarial attack during the training.

Although adversarial training is accepted as one of the most effective adversarial defense techniques that is applied to various domains [14], [15], it has received criticism due to several issues. First, the adversarial training conditions the robustness on only one type of attack [16]. Unfortunately, both single step and iterative adversarial generation techniques fail to find transferable attacks [17]. Moreover, adversarial training with a strong iterative attack does not necessarily correspond to improved robustness against various kinds of attacks [18]. Therefore, adversarial training suffers from overfitting certain first-order adversaries. Second, the adversarial training with strong attacks forces the model to capture certain features on extremely perturbed images. As a result, the test performance on natural images degrades severely as the attack strength in adversarial training increases. Lastly, the adversarial training techniques with first-order attacks are computationally expensive since extra back-propagation steps are essential to generate adversarial samples at each training step. Therefore, the quest for a more generalizable adversarial defense technique with less computational overhead and pristine natural accuracy abides.

To alleviate the overfitting issue, we propose synthesizing diverse adversarial perturbations during adversarial training. Generative models can be utilized to introduce diversity in the adversarial perturbations. Although generative methods are

Inci M. Baytaş is with the Department of Computer Engineering, Boğaziçi University, İstanbul, Turkey.

E-mail: inci.baytas@boun.edu.tr

Debayan Deb is with LENS Inc., Okemos, MI, USA.

E-mail: debayan@lenscorp.ai

extremely versatile in modeling distributions, robust models trained with generative attacks [19]–[21] has not reached the desired level of adversarial robustness and generalizability to natural samples compared to adversarial training with gradient-based attacks.

Instead, we propose a generative method that synthesizes one-step diverse perturbations without utilizing gradient of the classifier’s loss with respect to input image. The proposed generator is incorporated into a regularized adversarial end-to-end training framework for enhanced robustness to perturbations while maintaining superior generalizability to natural samples.

The contributions of the study are summarized below:

- A lightweight generator is designed to synthesize adversarial perturbations from a random vector. Without utilizing the input image, the generator can output diverse adversarial perturbations.
- The generator is trained to maximize the optimal transport distance between natural and synthesized adversarial samples. Therefore, the perturbation generation process does not depend on the class labels.
- The perturbation generator does not require the gradient of the classifier’s loss function with respect to the input image. In addition, the generator synthesizes adversarial attacks in a single step. Thus, the proposed robust training has less time complexity than adversarial training with iterative attacks.
- The objective of the proposed robust classifier is regularized with optimal transport distance between the natural and synthesized adversarial samples. Thus, the representation learning layers are guided to output features that follow similar distributions for natural images and their adversarial counterparts.
- The proposed end-to-end robust training approach does not sacrifice the natural accuracy while providing a robustness that is comparable with the state-of-the-art adversarial accuracy in literature.

The rest of this paper is organized as follows. In Section II, we overview state-of-the-art adversarial training methods with gradient-based and generative attacks. In Section III, the proposed framework is presented. The experimental results and analysis are discussed in Section IV. Final discussion is provided in Section V.

## II. RELATED WORK

The PGD attack has been one of the most effective white-box first-order attacks in literature [9]. As a result, PGD adversarial training is one of the most effective adversarial defense techniques. On the other hand, the PGD adversarial training has a poor generalization performance due to the lack of diversity in the PGD attacks [22]. There is a growing number of studies that investigate more robust and more generalizable adversarial defense techniques. While some of the studies aim to generate more diverse gradient-based attacks [23]–[25], some of them obtain adversarial samples through generative processes [19]–[21]. In this section, we overview studies that focus on improving the robustness and generalizability of adversarial training with gradient-based and generative attacks.

### A. Adversarial training with gradient-based attacks

Various studies in the literature have analyzed the trade-off between adversarial robustness and generalization. For instance, Zhang and Wang [23] discussed that the common issues in adversarial training, such as label leaking, are due to solely focusing on the decision boundary during attack generation. To facilitate a more generalizable adversarial training, adversarial perturbations are obtained considering inter-sample relationships rather than the decision boundary [23]. This feature scattering-based method generates adversarial samples via maximizing the difference between the distributions of natural and adversarial samples over a mini-batch. The feature scattering-based attack [23] enables improved robustness with a single-step gradient-based attack compared with the standard PGD adversarial training [9]. Wang and Zhang [26] also proposed bilateral adversarial training where both image and label are perturbed during training. The authors adopted targeted attack with the most confusing class.

Some studies argue that the existence of the adversarial samples is due to the features with non-robust components. Lee [24] proposed a vertex mixup approach to alleviate overfitting to the non-robust features. Their proposed adversarial vertex mixup approach comprises label smoothing and data augmentation. Label smoothing alleviates overfitting by making the model less confident about its predictions. Meanwhile, the training set is diversified through the combinations of adversarial and original samples, where the adversarial samples are obtained with PGD. Another interpolation technique to improve the adversarial robustness was studied by Zhang and Xu [25]. In their study, adversarial interpolation [25] is computed by minimizing the distance between the feature representations of a randomly perturbed sample and an adversarial sample, which is generated by maximizing the cross-entropy loss function. The interpolated input and target pairs are used in adversarial training. Adversarial interpolation is later used along with generative attacks to boost the generalization of the robust training [19].

Empirical analysis demonstrates that the generalization property of adversarial training can be improved without sacrificing robustness with methods such as Feature Scatter [23] and Adversarial Interpolation Training [25]. However, the challenges due to the complexity of incorporating a gradient-based attack and overfitting to a single type of attack remain. For this reason, the quest for more diverse and efficient attacks to integrate into the robust training leads the way to the generative models.

### B. Adversarial training with generative attacks

White-box gradient-based attacks have gained recognition as the most powerful adversarial perturbation generation technique. On the other hand, incorporating powerful adversarial attacks in robust training may not always reflect strong robustness. Since multi-step gradient attacks are prone to lack diversity, generative models have been considered an option for exploring adversarial samples in the input space.

Generative Adversarial Trainer (GAT) by Lee *et al.* [27] is one of the early defense techniques that replace the sign

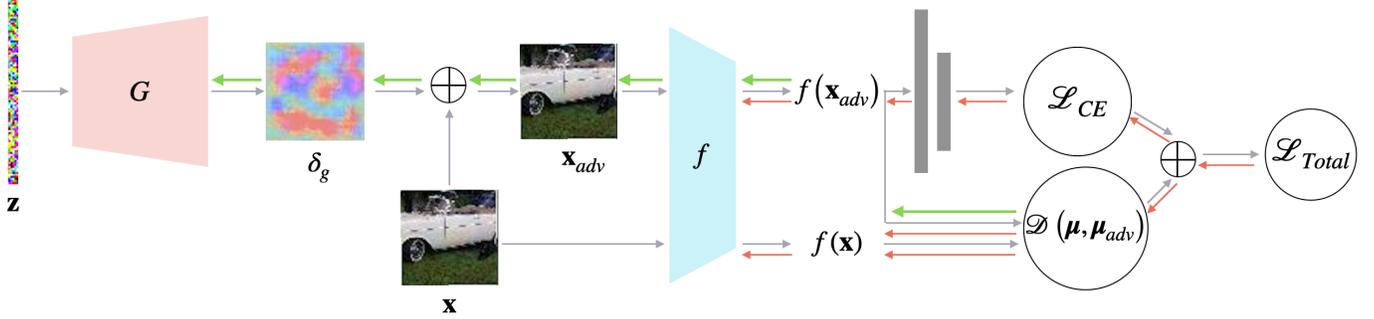


Fig. 1: Overview of the robustness via synthesis framework. Input of the generator is a normal random vector denoted by  $\mathbf{z}$ . The perturbation obtained by the generator  $G$  is denoted by  $\delta_g$ . The adversarial sample is computed by adding  $\delta_g$  to the original image  $\mathbf{x}$ . The latent representations of the original  $f(\mathbf{x})$  and the adversarial images  $f(\mathbf{x}_{adv})$  are obtained using a CNN architecture. The red and green arrows represent the path of error propagation to update the weights of  $f$  and  $G$ , respectively. The generator is updated using the optimal transport distance  $\mathcal{D}(\mu, \mu_{adv})$  between the latent representations of the original and the adversarial images. On the other hand, the classifier CNN is updated using the cross-entropy loss  $\mathcal{L}_{CE}$  regularized by  $\mathcal{D}(\mu, \mu_{adv})$ .

of the input gradient with a generator, which is conditioned on the input gradient. The amount of the perturbation is controlled by the  $\ell_2$  norm regularization. After the generator is updated, the training set is augmented with the adversarial samples. The discriminator’s objective function is the same as the objective function with adversarial regularization proposed by Goodfellow *et al.* [13]. Samangouei *et al.* [28], on the other hand, posed the robustness against adversarial attacks as denoising adversarial samples encountered during inference. Therefore, the authors designed Defense-GAN, which can be used with any classifier. At test time, the Defense-GAN generates a reconstruction of the original input without the adversarial noise. However, the empirical evidence for the denoising capabilities of the Defense-GAN is only demonstrated for one-step FGSM and CW attacks. For this reason, the efficiency of the Defense-GAN against multi-step attacks is unclear.

One of the factors that cause poor generalization is using one type of attack to solve the inner maximization problem of the adversarial training. Dong *et al.* [29] addressed this issue by learning an adversarial distribution rather than a single adversarial sample to solve the inner maximization problem. Their proposed solution, namely Adversarial Distribution Training (ADT), models the distribution of the potential adversarial samples around each input. Compared with the state-of-the-art results on benchmark datasets, the performance of the ADT is inferior. However, the level of robustness provided by ADT is more consistent across various attacks.

In a more recent study, Jeddi *et al.* [21] proposed a robust training framework, namely Learn2Perturb, by injecting perturbations into each layer to increase the uncertainty of the deep network. The perturbations, which are added to each feature map, are drawn from a zero-mean normal distribution with a learnable standard deviation. Learn2Perturb alternately updates the parameters of the model and perturbation-injection. Unlike the traditional generation of adversarial perturbations, the Learn2Perturb aims to learn the distribution of the perturbations in the latent spaces governed by the feature

maps. On the other hand, the Learning2Perturb model cannot overcome the poor generalization performance of the PGD adversarial training considering its test performance on natural images.

Jiang *et al.* [19] on the other hand, adopted a generic learning-to-learn (L2L) framework for adversarial training. The authors designed an attacker network to generate perturbations. The attacker network either synthesizes perturbations via only the original samples or concatenation with their gradients. The latter is called gradient attacker, and its multi-step version is also provided [19]. The multi-step gradient attacker aims to mimic the PGD attack with an RNN. Although one-step and two-step gradient attackers improve the adversarial accuracy, they still cannot generalize well to the natural samples. However, Jiang *et al.* [19] also showed that combining L2L with adversarial interpolation training yields a boost in both natural and adversarial accuracies.

L2L based robust training is also investigated by Jang *et al.* [20]. The authors proposed a generator that can synthesize strong and diverse attacks. Their proposed framework, named L2L-DA [20], adopts a recursive approach to generate strong perturbations while the diversity is enforced by an additional diversity loss. Unlike the gradient attacker in L2L [19], the generator takes random noise in addition to the original sample and its gradient. Experiments on benchmark datasets indicate improved robustness compared to L2L [19]. However, both robust and natural accuracies of L2L-DA cannot reach state-of-the-art performances.

Wang and Yu [30] also designed a GAN to parametrize the inner maximization problem. A more traditional GAN structure is considered to generate adversarial perturbations. Similar to L2L-DA [20] and L2L [19], the input of the generator is the original sample. The  $\tanh$  activation function in the last layer of the generator ensures that the perturbations will not exceed the epsilon ball. The author also suggested regularizing the discriminator network with the norm of the gradient in order to stabilize the GAN training. Their proposed robust training cannot improve the adversarial accuracy over

PGD adversarial training, however, it achieves a much higher natural accuracy than the PGD adversarial training.

The majority of the studies reviewed in this section require the gradient of the classifier’s loss. Even the generative perturbation frameworks take gradient information along with the input sample and utilize recurrent architectures to intensify the adversarial perturbations. However, to the best of our knowledge, robust training techniques with generative frameworks sacrifice either robustness or generalization to natural samples in order to increase the attack strength and diversity. In this study, we present an adversarial perturbation generation method that promotes diversity without utilizing any input sample or gradient information. Yet, the proposed framework provides a more robust and generalizable classifier than the adversarial training techniques with generative attacks discussed in this section. In this regard, we posit that regularizing the classifier is crucial to lead the generator to explore diverse adversarial directions. Updating the classifier with cross-entropy loss alone is not sufficient to learn a robust representation even if we augment the training set with strong generative attacks.

### III. METHOD

Due to the fact that we do not have access to all possible samples that lie in the input space (including adversarial examples), the standard training procedure inevitably encodes non-robust features. Therefore, if we require any robustness against adversarial directions in the input space, it is imperative to introduce diverse perturbations to the classifier during training. To promote diversity in perturbations, we propose a regularized adversarial training technique that encourages the proposed generator to synthesize diverse attacks.

#### A. Problem Definition

Given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  of  $N$  natural images  $\mathbf{x}_i \in \mathbb{R}^{d_w \times d_h \times c_{in}}$  and their labels  $y_i \in 1, \dots, C$ , the adversarial training is posed as the following two step optimization problem:

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \max_{\delta_i \in S} \mathcal{L}(f(\mathbf{x}_i + \delta_i; \theta), y_i) \quad (1)$$

where  $\theta$  is the parameter of the classifier network,  $\delta_i$  is the perturbation generated for  $\mathbf{x}_i$ ,  $S = \{\delta : \|\delta\|_{\infty} \leq \epsilon\}$  is the set of allowed perturbations, and  $\mathcal{L}(\cdot)$  is an objective function, e.g, cross-entropy.

The perturbation  $\delta$  is generated via the inner maximization step in the Eq. 1, which is an intractable problem. The gradient of the objective function at an input data point is extremely informative about the adversarial direction. For this reason, one of the most powerful adversarial attack approaches is Projected Gradient Descent (PGD) [9] given below.

$$\mathbf{x}^{t+1} = \Pi_{\mathbf{x}+S}(\mathbf{x}^t + \alpha \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f(\mathbf{x}_i + \delta_i; \theta), y_i))) \quad (2)$$

where  $\Pi$  is a projection ensuring that the amount of perturbation will not exceed  $\epsilon$  and the adversarial sample will be inside the input domain,  $\alpha$  is the step size, and

$\text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f(\mathbf{x}_i + \delta_i; \theta), y_i))$  is the direction at the input that increases the value of the objective function. Thus, the PGD adversarial training [9] poses robust training as a saddle point problem, where the PGD attacks are obtained by the maximization step and the model parameters are updated via the minimization problem.

There are several issues with adversarial training when the attack method is PGD. Although the PGD attack is quite strong, the diversity among the attacks at different training iterations is limited. Thus, label leaking and catastrophic overfitting often hinder both the robust and generalization performance of the PGD adversarial training. Besides, when the perturbations obtained by PGD are added to the input image, we may lose the information of the natural patterns. We can observe this issue by the degradation in the natural accuracy of the adversarially trained models. Furthermore, iterative gradient-based attacks are very time-consuming due to multiple backpropagation steps at each iteration.

In this study, the proposed adversarial training framework, illustrated in Fig. 1, is designed to induce adversarial robustness without sacrificing the natural accuracy. In particular, the proposed approach aims to alleviate the following challenges.

- Although the iterative gradient-based adversarial attacks [9] are strong regarding the reduction in test accuracy, they often push the natural samples towards similar adversarial directions such that the perturbations may be constantly projected back onto the boundary of  $\ell_{\infty}$  norm ball. Thus, the adversarial perturbations are not versatile enough to enable robustness without overfitting to a certain type of attack.
- Iterative gradient-based adversarial attack generation algorithms are computationally expensive. Therefore, they may not scale to large-scale problems.
- In robust training with generative attacks, when the attack generator network is conditioned on a single input or its gradient, the generator network may overfit to a certain subtle perturbation around the input sample. This situation gets more acute when the parameters of the generator are updated by maximizing the cross-entropy loss.
- Updating both attack generator and the classifier networks with cross-entropy loss *alone* results in weaker perturbations.

The proposed framework comprises two modules: adversarial perturbation generator and a classifier. In the next section, adversarial perturbation generator will be discussed.

#### B. Adversarial Perturbation Generator

Generative Adversarial Networks (GANs) have been one of the most popular generative models that indirectly learn the distribution of the input data. The generator network may take a random vector or may be conditioned on a specific sample and generates realistic fake samples. A well-trained binary classifier that discriminates a sample as fake or real, is necessary to induce the generator to output more realistic samples [31]. Inspired by the GAN mechanism, we propose

TABLE I: Generator architecture. The latent representation before the last fully connected layer of the WRN-28-10 network is 640 dimensional. The CIFAR10, CIFAR100, and SVHN datasets have RGB images of size  $32 \times 32 \times 3$ . For this reason, a normal random vector is drawn from the 640 dimensional space and decoded to a  $32 \times 32 \times 3$  tensor.

Layer	Size	Output
Input	$\mathbf{z} \in \mathbb{R}^{640} \sim \mathcal{N}(0, 1)$	$640 \times 1$
FC	$640 \times 4096$	$4096 \times 1$
Reshape		$8 \times 8 \times 64$
DeConv	Kernel:4 × 4, Stride:2, Channels:32	$16 \times 16 \times 32$
Batchnorm		
Leaky ReLU	Leakiness:0.2	
DeConv	Kernel:4 × 4, Stride:2, Channels:16	$32 \times 32 \times 16$
Batchnorm		
Leaky ReLU	Leakiness:0.2	
Conv	Kernel:4 × 4, Stride:1, Channels:3	$32 \times 32 \times 3$

an adversarial perturbation generator  $G(\mathbf{z}; \Phi)$  where  $\mathbf{z}$  is a random vector and  $\Phi$  is the parameters of the generator.

The main purpose of using a generator to obtain adversarial perturbations is to learn how the perturbations within an  $\epsilon$  ball are distributed. As seen in Fig. 1, unlike many generative attack models in literature, the proposed generator  $G(\mathbf{z}; \Phi)$  is not conditioned on the input sample or the gradient of the loss function at the input. There are two essential reasons behind this design. First, the proposed generator does not output an adversarial image but an adversarial perturbation. For this reason, the output of the generator is not explicitly forced to be visually similar to a particular pattern. Consequently, generating the perturbation from a random vector facilitates diversity among adversarial samples compared with generating a perturbation at a particular data point. Second, decoding a random vector into a perturbation tensor is less computationally expensive than including an extra encoder in a pixel to pixel setting.

To generate an adversarial attack, we first sample a normal random vector,  $\mathbf{z} \in \mathbb{R}^d$ , in the latent space of the classifier. The random vector is then decoded into a tensor of the same size as the input image  $\mathbf{x}$  as shown in Table I. This tensor has unbounded values. However, we avoid using the tanh activation function in the last layer of the generator. Due to its flat regions, the tanh function is notoriously prone to numerical instabilities. Moreover, according to our observations, the tanh activation function forces the majority of the perturbation values to be exactly  $-1$  and  $1$  that prevents the diversity. To ensure that the perturbation is within the set of allowed perturbations in the  $\ell_\infty$  ball, clipping between the values  $[-\epsilon, \epsilon]$  is employed instead of multiplying the tanh output by  $\epsilon$ . Thus, the generator is intended to be more flexible to explore potential adversaries in a wider territory.

After obtaining the perturbation  $\delta$ , the adversarial sample is obtained as follows

$$\delta_g = \Pi_S(G(\mathbf{z}; \Phi)) \quad (3)$$

$$\mathbf{x}_{\text{adv}} = \Pi_{\mathbf{x}}(\mathbf{x} + \delta_g) \quad (4)$$

where  $\mathbf{x}_{\text{adv}} \in \mathbb{R}^{d_w \times d_h \times c_{in}}$  is the adversarial sample,  $\Pi_S$  and  $\Pi_{\mathbf{x}}$  denote the projection operators (e.g., clipping) to map the

perturbation into  $S$  and the adversarial sample back into the input domain, respectively.

Parameters of the generator,  $\Phi$  are updated through the following optimization problem.

$$\max_{\Phi} \mathcal{D}(\mu, \mu_{\text{adv}}) \quad (5)$$

where  $\mathcal{D}(\mu, \mu_{\text{adv}})$  denotes the optimal transport (OT) distance between the natural and the adversarial sample distributions. The OT distance is a well-studied distance that stabilizes and improves the GAN training [32], [33]. The OT distance between two distributions is defined as follows.

$$\mathcal{D}(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} E_{(\mathbf{x}, \mathbf{y}) \sim \gamma} c(\mathbf{x}, \mathbf{y}) \quad (6)$$

where  $\Pi(\mu, \nu)$  is the set of joint distributions  $\gamma(\mathbf{x}, \mathbf{y})$  with marginal distributions of  $\mu(\mathbf{x})$  and  $\nu(\mathbf{y})$ , and  $c(\mathbf{x}, \mathbf{y})$  is a cost function [23]. The OT distance represents the minimum cost to transport one marginal to another.

Following the footsteps of Zhang *et al.* [23], finding the OT distance is equivalent to the following problem.

$$\mathcal{D}(\mu, \mu_{\text{adv}}) = \min_{\mathbf{T} \in \Pi(\mathbf{u}, \mathbf{u}_{\text{adv}})} \sum_{i=1}^n \sum_{j=1}^n \mathbf{T}_{ij} \cdot c\left(f(\mathbf{x})_i, f(\mathbf{x}_{\text{adv}})_j\right) \quad (7)$$

where  $n$  denotes the size of the mini-batch,  $\Pi(\mathbf{u}, \mathbf{u}_{\text{adv}}) = \{\mathbf{T} \in \mathcal{R}_+^{n \times n} | \mathbf{T} \mathbf{1}_n = \mathbf{u}, \mathbf{T}^T \mathbf{1}_n = \mathbf{u}_{\text{adv}}\}$ ,  $\mathbf{1}_n$  is  $n$ -dimensional all-one vector,  $\mathbf{u}$  and  $\mathbf{u}_{\text{adv}}$  contain the values of the weight vectors  $\mu = \{u_i\}_{i=1}^n$  and  $\mu_{\text{adv}} = \{v_i\}_{i=1}^n$ , respectively [23]. Since  $\mu$  and  $\mu_{\text{adv}}$  are probability distributions,  $\sum_{i=1}^n u_i = \sum_{i=1}^n v_i = 1$ . In this study, the cost function in Eq. 7 is defined as the euclidean distance between the latent representations of the natural and the adversarial samples as shown below.

$$c\left(f(\mathbf{x})_i, f(\mathbf{x}_{\text{adv}})_j\right) = \|f(\mathbf{x})_i - f(\mathbf{x}_{\text{adv}})_j\|_2^2 \quad (8)$$

The OT distance in Eq. 7 is coupled over mini-batches such that the distance between empirical distributions of the natural samples and their adversarial counterparts in one mini-batch is measured. Since the adversarial perturbation is not generated at a single sample, the generation process potentially explores a wider region in the input space. Secondly, the perturbation is not generated by taking the sign of the gradient of the OT distance. When we add the sign of the gradient to the input sample, every pixel moves exactly the same amount. For the pixel values that are close to the upper and lower bounds of the input domain, the perturbation may not have any effect due to clipping. Such cases are prone to overfitting. Therefore, a generator is designed to output a perturbation tensor whose values range between  $[-\epsilon, \epsilon]$ .

The proposed generator does not require multiple backpropagation steps or a recurrent loop. Thus, the proposed adversarial perturbation generator is less expensive than iterative attacks. As presented in Algorithm 1, at each iteration, the perturbation generator is updated once. The generator weights are updated in order to maximize the OT distance between the latent representations of natural and synthesized adversarial samples. As the classifier becomes more robust during the training, the perturbation generator will be encouraged to

explore more sophisticated perturbations. In the next section, details of the classifier training are presented.

### C. Classifier

In a standard GAN setting, a discriminator is trained to classify fake and real images while providing informative gradients to the generator. In this study, unlike the standard-setting, a classifier is trained to accomplish an object recognition task. Inspired by the GAN pipeline, the OT distance loss is backpropagated through the classifier to update the generator’s parameters. For this reason, the status of the classifier during the training enables the generator to create diverse adversarial perturbations. Furthermore, the goal of robust training is to learn a latent space in which the feature representations of the natural sample and its adversarial counterpart do not differ. However, updating the classifier’s parameters based solely on the cross-entropy loss may not achieve this objective.

To alleviate the challenge discussed above, the OT distance between the natural and the synthesized adversarial sample distributions is used as a regularizer. The total loss function of the classifier is given below.

$$\mathcal{L}_{\text{total}}(\mathbf{x}, \mathbf{x}_{\text{adv}}, y; \theta, \Phi) = \mathcal{L}_{CE}(\mathbf{x}_{\text{adv}}, y; \theta) + \mathcal{D}(\mu, \mu_{\text{adv}}) \quad (9)$$

where  $\mathcal{L}_{CE}(\mathbf{x}, y; \theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{ij} \log f_j(\mathbf{x}_i; \theta)$  is the cross-entropy loss function,  $\theta$  denotes the parameters of the classifier, and  $\mathcal{D}(\mu, \mu_{\text{adv}})$  is defined in Eq. 7.

A significant improvement in robustness is observed when the OT distance regularization is included in the loss function of the classifier given in Eq. 9. The contributions of the OT distance regularization is two-fold; (i) classifier is encouraged to learn a latent space that is less susceptible to adversarial perturbations, (ii) generator is encouraged to produce more challenging perturbations since the classifier is trained to map adversarial samples near to their natural counterparts in the latent space.

In the robust training setting, feeding the classifier with very strong adversarial attacks such as PGD hurts the generalization performance and overfits to the particular attack method. For this reason, adversarial training demands exploring diverse adversarial samples in the input space rather than samples on the boundary of the  $\epsilon$ -ball that maximizes the objective function. In the next section, the proposed robustness via synthesis framework given in Algorithm 1 is analyzed with experiments on different benchmark datasets.

## IV. EXPERIMENTS

The robustness of the proposed approach is tested on three common benchmark datasets, namely CIFAR10, CIFAR100 [34] and SVHN [35]. In our experiments, object recognition task is considered and classification accuracy is used to compare the proposed and baseline methods. In addition to the white-box and black-box performances, diversity in the perturbation generation and the behavior of the learned latent representations by the proposed model are analyzed in Section IV-C3. We also investigate the effect of OT loss regularization and the choice of the loss function to train the generator with an ablation study in Section IV-C4.

---

**Algorithm 1** Robustness via synthesis for T epochs, M mini-batches,  $\theta$  network parameter,  $\Phi$  generator parameter, and  $\mathcal{L}_{\text{total}}$  loss function

---

```

1: for  $t = 1, \dots, T$  do
2:   for  $i = 1, \dots, M$  do
3:     Sample a normal random vector,  $\mathbf{z}$ 
4:      $\delta_g = \Pi_S(G(\mathbf{z}; \Phi))$ 
5:      $\mathbf{x}_{\text{adv}} = \Pi_{\mathbf{x}}(\mathbf{x} + \delta_g)$ 
6:      $\Phi^* = \operatorname{argmax}_{\Phi} \mathcal{D}(\mu, \mu_{\text{adv}})$ 
7:      $\delta_g = \Pi_S(G(\mathbf{z}; \Phi^*))$ 
8:      $\mathbf{x}_{\text{adv}} = \Pi_{\mathbf{x}}(\mathbf{x} + \delta_g)$ 
9:      $\theta^* = \operatorname{argmin}_{\theta} \mathcal{L}_{\text{total}}(\mathbf{x}, \mathbf{x}_{\text{adv}}, y; \theta, \Phi^*)$ 
10:   end for
11: end for

```

---

### A. Datasets and Implementation Details

We conduct experiments with CIFAR10, CIFAR100 [34], and SVHN [35] benchmark datasets of  $32 \times 32$  RGB images. The CIFAR10 and SVHN datasets have 10 classes, whereas the CIFAR100 dataset has 100 categories. The CIFAR10 and CIFAR100 datasets have 50k training images and 10k test images. The SVHN dataset has 73,257 training images and 26,032 test images. In our experiments, compatible with the adversarial training literature, WRN-28-10 wide ResNet [37] classifier is used. The proposed end-to-end robustness-via-synthesis training framework is implemented in Tensorflow within the codebase provided by MadryLab [38]. Tensorflow implementation of the Sinkhorn algorithm with regularization of 0.01 and 100 iterations [39] is used to compute the OT distance.

For all of the experiments, the maximum perturbation amount is set to  $\epsilon = 8/255$ . For the CIFAR datasets, random cropping and flipping are applied to augment the training sets [38]. Throughout the experiments, the batch size is fixed to 64. The CIFAR10 and SVHN models are trained for 179k iterations at which the most optimal robustness is observed. The CIFAR100 model is trained for 300k iterations. Since the OT distance is computed over mini-batches, datasets with a large number of categories may require longer training in order to present all the categories to the generator.

As suggested in literature [23], [26], a learning rate scheduling scheme is adopted with transition steps of 60K and 90K iterations. An initial learning rate of 0.1 is used for the CIFAR datasets and 0.01 for the SVHN dataset with a learning rate decay of 0.1. While updating the parameters of the discriminator, label smoothing with a weight factor of 0.5 is applied. The learning rate for the proposed generator architecture, provided in Table I, is set to 0.01 and 0.0001 for CIFAR datasets, and the SVHN dataset, respectively. A learning rate schedule is not applied during the optimization of the generator. The momentum optimizer is utilized to train both the classifier and the generator.

### B. Baselines

The proposed robustness via synthesis technique is an adversarial training approach with a generative attack. We compare

TABLE II: White-box Results. In all of the experiments, the maximum amount of perturbation is set to  $\epsilon = 8/255$ , and the step size is  $2/255$ . Starred values are obtained by training the corresponding networks from scratch. Other values represent the best performances reported by the baseline studies. The values of the cells with a dash are not available. CIFAR100 results of Learn2Perturb are marked with † to indicate that these values are taken from Figures 5 and 6 in the supplementary material of the study [21].

CIFAR10									
Defenses	Natural	FGSM	PGD-7	PGD-10	PGD-20	PGD-100	CW-20	CW-100	AdvGAN [36]
Natural	94.42*	10.22*	0	0	0	0	0	0	10.89*
AT [9]	87.25	62.64	49.67	47.33	45.91	45.29	46.99	46.54	<b>84.98*</b>
Bilateral [26]	91.00	70.70	63.00	-	57.80	55.20	56.20	53.80	-
FeatScatter [23]	90.00	78.40	<b>73.54*</b>	<b>70.90</b>	70.50	68.60	62.40	60.60	-
Adv-interp [25]	90.30	78.00	-	-	<b>73.50</b>	<b>73.00</b>	<b>69.70</b>	<b>68.70</b>	-
Adv-vertex [24]	93.24	78.25	-	62.67	58.23	-	53.63	-	-
L2L-DA [20]	78.91	45.77	-	39.69	-	38.39	-	37.75	-
L2L [19]	85.35	-	-	-	54.32	52.12	-	57.07	-
Direct [30]	91.08	72.81	-	44.28	-	-	-	-	-
Learn2Perturb [21]	85.30	62.43	56.06	-	-	-	-	-	-
<b>Ours</b>	<b>94.17</b>	<b>79.99</b>	69.02	65.24	57.50	48.62	41.21	27.54	72.44

CIFAR100									
Defenses	Natural	FGSM	PGD-7	PGD-10	PGD-20	PGD-100	CW-20	CW-100	AdvGAN [36]
Natural	79.22*	3.52*	0	0	0	0	0	0	4.14*
AT [9]	59.78*	32.70*	25.07*	23.49*	22.78*	22.44*	23.05*	22.87*	<b>55.39*</b>
Bilateral [26]	66.20	31.30	-	-	23.10	22.40	-	20.00	-
FeatScatter [23]	73.90	61.00	46.29*	45.99*	47.20	<b>46.20</b>	34.60	30.60	-
Adv-interp [25]	73.6	58.3	-	-	41	40.2	32.4	31.2	-
Adv-vertex [24]	74.81	<b>62.76</b>	-	-	38.49	-	-	-	-
L2L [19]	60.95	-	-	-	31.03	29.75	-	32.28	-
Direct [30]	70.99	41.86	-	18.25	-	-	-	-	-
Learn2Perturb [21]	58.00†	30.00†	26.00†	-	-	-	-	-	-
<b>Ours</b>	<b>76.32</b>	51.63	49.50	<b>49.15</b>	<b>48.16</b>	45.79	<b>40.60</b>	<b>38.25</b>	54.52

SVHN									
Defenses	Natural	FGSM	PGD-7	PGD-10	PGD-20	PGD-100	CW-20	CW-100	AdvGAN [36]
Natural	95.92*	13.62*	0	0.18*	0	0	0	0	45.30*
AT [9]	90.74*	64.50*	47.87*	44.23*	41.38*	40.37*	42.46*	41.60*	88.10*
Bilateral [26]	94.10	69.80	-	-	53.90	50.30	-	48.90	-
FeatScatter [23]	96.20	83.50	61.88*	55.40*	62.90	52.00	61.30	50.80	-
Adv-interp [25]	94.10	75.60	-	-	<b>65.80</b>	<b>64.00</b>	<b>63.40</b>	<b>60.40</b>	-
Adv-vertex [24]	95.59	81.83	-	-	61.90	-	-	-	-
Direct [30] ( $\epsilon = 0.05$ )	<b>96.34</b>	<b>91.51</b>	-	37.97	-	-	-	-	-
<b>Ours</b>	95.50	80.37	70.17	<b>67.13</b>	60.53	53.98	54.21	46.33	82.67

our method to the recent adversarial training techniques that are in the same category as the proposed method, such as L2L-DA [20], L2L [19], Learn2Perturb [21], and Direct [30]. In Table II, the best performances reported by the studies [19]–[21], [30] are included.

Although the proposed framework is designed to improve the performance of adversarial training with generative attacks, we also provide a comparison with the state-of-the-art adversarial training techniques with gradient-based attacks, such as standard adversarial training (AT) proposed by Madry *et al.* [9], TRADES [18], bilateral adversarial training (Bilateral) [26], feature scatter-based adversarial training (FeatScatter) [23], and adversarial interpolation training (Adv-interp) [25]. For CIFAR100 and SVHN, we train the AT [9] model from scratch. In AT [9] experiments, implementation and hyperparameters provided by MadryLab [38] are utilized. During the adversarial training, PGD-7 attack is utilized. For FeatScatter [23] experiments, we trained the models from scratch using the Pytorch implementation [40] for CIFAR100 and SVHN to report the PGD-10 results. However, we present the best results reported by the authors when they are available.

### C. Classification Performance

In this section, we compare the classification accuracy under various types of attacks. In Table II, all the attacks are generated given the architecture and the parameters of the models. In Table III, adversarial attacks are generated using surrogate models and tested on robust models.

1) *White-Box Attack*: PGD [9] is considered as one of the most effective first-order white-box adversarial attacks. Robustness performance of the proposed and the baseline methods are evaluated against FGSM [13], PGD-10, PGD-20, and PGD-100 for  $\ell_\infty$  attacks, and CW-20 and CW-100 [4] for the  $\ell_2$  norm attacks. The numbers next to the attack types indicate the number of steps. Random initialization is applied to all of the white-box attacks and the step size is set to  $2/255$ .

In addition to the gradient-based attacks, white-box performance is also evaluated against a prominent generative attack, namely AdvGAN [36]. To obtain the adversarial samples via AdvGAN, the generator of the AdvGAN is trained from scratch to fool robust classifiers obtained by the proposed framework, AT [9], and the natural model. We used the

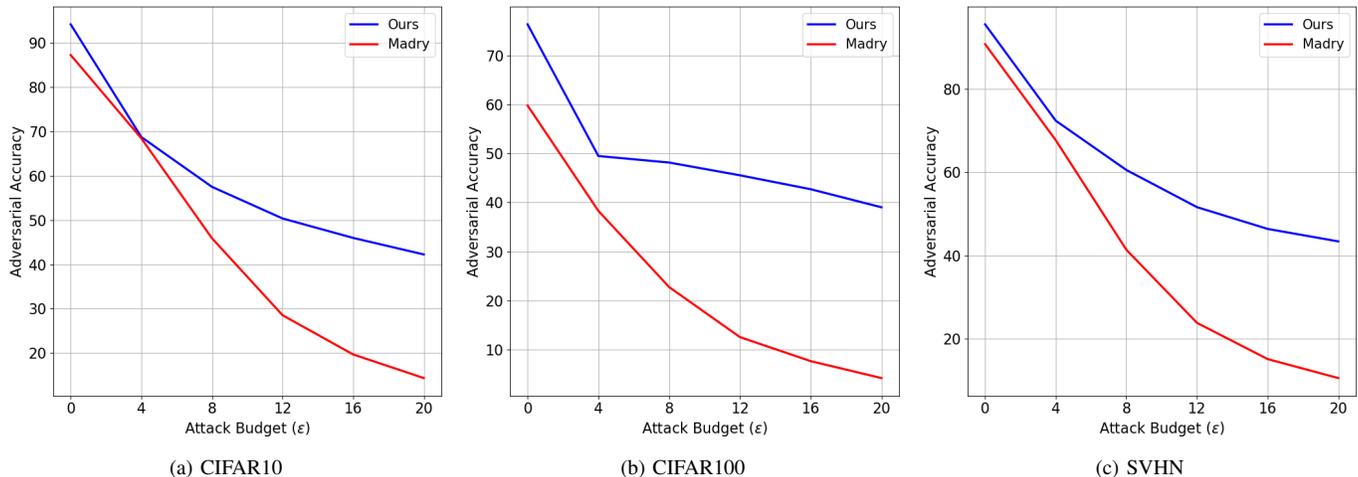


Fig. 2: Robustness across varying attack budgets for (a) CIFAR10, (b) CIFAR100, and (c) SVHN. Adversarial accuracies against the PGD-20 attack of the proposed and Madry’s models [9] are plotted for varying  $\epsilon \in [\frac{0}{255}, \frac{20}{255}]$  values. The step size of the PGD attack is set to  $\frac{2}{255}$ .

Tensorflow implementation of AdvGAN<sup>1</sup>. In our experiments, we train the generator and the discriminator of the AdvGAN from scratch and report the result at the checkpoint where the adversarial attack is the most successful.

As seen in the white-box results, the proposed approach can preserve the natural accuracy across all the datasets. Particularly for CIFAR10 and CIFAR100, generative methods [19], [20], [30] fail to generalize to the natural samples and their robustness is poor compared to gradient-based methods. On the other hand, the robustness obtained via the proposed framework is not only better than the other generative approaches and Madry’s AT [9], but also comparable to the state-of-the-art performances by gradient-based techniques such as Feature Scatter [23]. Reaching this level of robustness and generalizability without the gradient of the cross-entropy loss in the perturbation generation is a notable point. This can be clearly observed when we investigate the change in robustness against PGD-20 for varying attack budgets in Figure 2. We find that the proposed framework can maintain much better robustness than the standard AT even for high  $\epsilon$  values. For instance, the proposed framework facilitates preserving the robustness for  $\epsilon = 20/255$ .

2) *Black-Box*: Black-box performance is evaluated against one generative and several gradient-based attacks. The gradient-based black-box attacks are generated using a naturally trained model and robust model that is obtained with the standard adversarial training with PGD-7 [9]. Using each model, a set of PGD and CW attacks with step sizes of 20 and 100 were generated. Similarly, we generate adversarial samples via AdvGAN using the natural and robust models to evaluate the black-box performance of the proposed model. As seen in Table III, the proposed framework is generalizable to various black-box scenarios.

3) *Perturbation Diversity Analysis*: To improve the generalizability of the robust model, diverse adversarial samples

TABLE III: Black-box Results. The Natural Model table presents the robustness of the proposed model against adversarial samples that are generated using the pre-trained natural model. Whereas, Robust Model table contains the performance of the proposed model against adversarial samples generated using the robust model of Madry *et al.* [9].

Natural Model			
Attacks \ Datasets	CIFAR10	CIFAR100	SVHN
PGD-20	71.11	62.10	82.42
PGD-100	71.46	61.51	83.38
CW-20	71.57	63.57	82.54
CW-100	71.76	63.77	82.63
AdvGAN [36]	85.08	52.19	95.72

Robust Model			
Attacks \ Datasets	CIFAR10	CIFAR100	SVHN
PGD-20	74.53	54.34	65.69
PGD-100	74.12	54.55	65.16
CW-20	75.31	55.01	66.69
CW-100	74.96	55.23	66.46
AdvGAN [36]	88.22	68.10	92.58

should be used in the training. One of the motivations for generative perturbations is to explore more diverse adversarial attacks than gradient-based techniques. In this section, we investigate the diversity of the adversarial perturbations generated during the proposed regularized robust training method in the latent space. For this purpose, we generate 100 adversarial images for each natural sample in the CIFAR10 test set. The adversarial images are obtained from 100 different random vectors. Then, the Euclidean distances between the normalized latent representations of the natural sample and 100 adversarial samples are computed. In Figure 3, we plot the average distances over all the data points in the test set for 174 checkpoints. The curve represents the change in the average distance during training. The shaded region around the curve represents the average standard deviation of the distances

<sup>1</sup><https://github.com/ctargon/AdvGAN-tf>

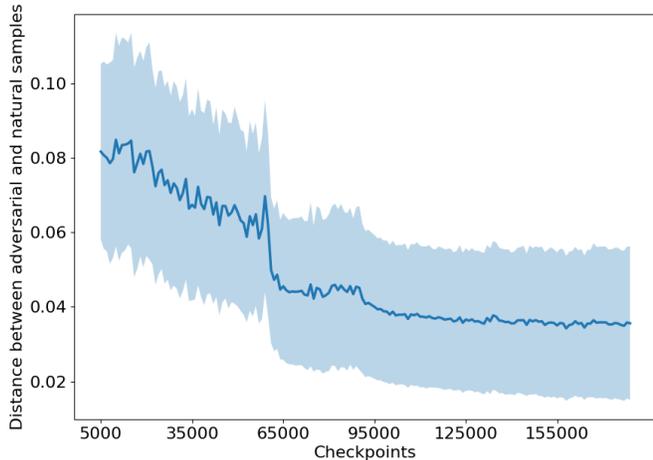


Fig. 3: We generate 100 adversarial images for each natural sample in the CIFAR10 test set. The adversarial images are obtained from 100 different random vectors. Then, the Euclidean distances between the normalized latent representations of the natural sample and 100 adversarial samples are computed. The average distances for 174 checkpoints are plotted. The curve represents the change in the average distance during training. The shaded region around the curve represents the average standard deviation of the distances between the natural sample and its 100 adversarial samples. While the diversity in the adversarial directions is preserved during the robust training, the average distance between the latent representations decreases as expected.

between the natural sample and its 100 adversarial samples.

According to our initial hypothesis, there are two expected behaviors in this experiment: (i) The diversity in the adversarial perturbation generation should be preserved during the robust training. (ii) As the robust training progresses, the distances between adversarial samples and the natural samples should decrease since the learned representations are expected to be unsusceptible to adversarial perturbations. As it can be observed in Figure 3, the average distances decrease during the training while the standard deviation of the distances between the natural image and its adversarial samples is preserved.

We also demonstrated the behavior of the latent space in Figure 4. In this figure, we randomly sample a data point from each category in the test set of CIFAR10. Similar to the above experiment, we generate 100 different adversarial samples corresponding to the natural image that is randomly chosen. Then, the latent representations of the natural and their adversarial counterparts are investigated using the t-SNE plot in Figure 4. In the early stages of the robust training, the domain gap between the natural and adversarial samples is visible in the latent space. As the model becomes more robust, it is harder for the generator to find adversarial directions. Thus, for the majority of the categories, the adversarial samples are grouped such that their center is the natural sample.

4) *Ablation Study*: In this section, we investigate the necessity of the OT distance regularization and updating the weights

TABLE IV: Ablation Results. **noReg + OT** denotes the model without OT distance regularization such that the generator is updated with OT distance between the natural and adversarial samples and the classifier is updated with only the cross-entropy loss. **OT-Reg + Xent** denotes the model with OT distance regularization where the generator is updated with the cross-entropy loss instead of the OT distance. Finally, **OT-Reg + OT** denotes the proposed model with OT distance regularization and the generator being updated with OT distance.

CIFAR10			
Defenses	Natural	PGD-20	CW-20
<b>noReg + OT</b>	94.46	24.55	21.09
<b>OT-Reg + Xent</b>	94.44	52.01	36.70
<b>OT-Reg + OT (Proposed)</b>	94.17	57.50	41.21

CIFAR100			
Defenses	Natural	PGD-20	CW-20
<b>noReg + OT</b>	77.40	26.33	8.98
<b>OT-Reg + Xent</b>	76.95	42.28	30.32
<b>OT-Reg + OT (Proposed)</b>	76.32	48.16	40.60

SVHN			
Defenses	Natural	PGD-20	CW-20
<b>noReg + OT</b>	96.06	6.66	7.06
<b>OT-Reg + Xent</b>	96.26	29.29	24.14
<b>OT-Reg + OT (Proposed)</b>	95.50	60.53	54.21

of the generator with the OT distance rather than cross-entropy. In Table IV, **noReg + OT** denotes models trained without OT distance regularization, and **OT-Reg + Xent** represents models with the OT distance regularization in which the generator is updated with the cross-entropy loss. As can be observed in the table, OT distance regularization significantly improves the robustness of the model across all the datasets. In the presence of the regularizer, if the generator is updated with the cross-entropy loss, the adversarial accuracies of all three datasets decrease. This decrease is substantial for the SVHN dataset, which has completely different patterns than CIFAR10 datasets.

## V. CONCLUSION

In this study, a robust training approach with generative adversarial perturbations is proposed. The attack generation does not require pixel to pixel translation or recurrent architectures. More importantly, the proposed generator does not take any gradient information as input. Thus, the computational complexity is reduced compared to the robust training models with iterative gradient computations. To encourage diverse adversarial perturbations during training, the attack generation network is updated by maximizing the optimal transport distance between the representations of the synthesized adversarial and natural samples. Furthermore, the perturbation is generated from a random vector. As a result, the dependency on a single sample and its label is reduced during the attack generation. The optimal transport distance between the adversarial and natural samples is also utilized to regularize the classifier such that the learned representations are encouraged to be robust against adversarial perturbations. Experiments on CIFAR10, CIFAR100, and SVHN datasets demonstrated that the proposed robust training approach can

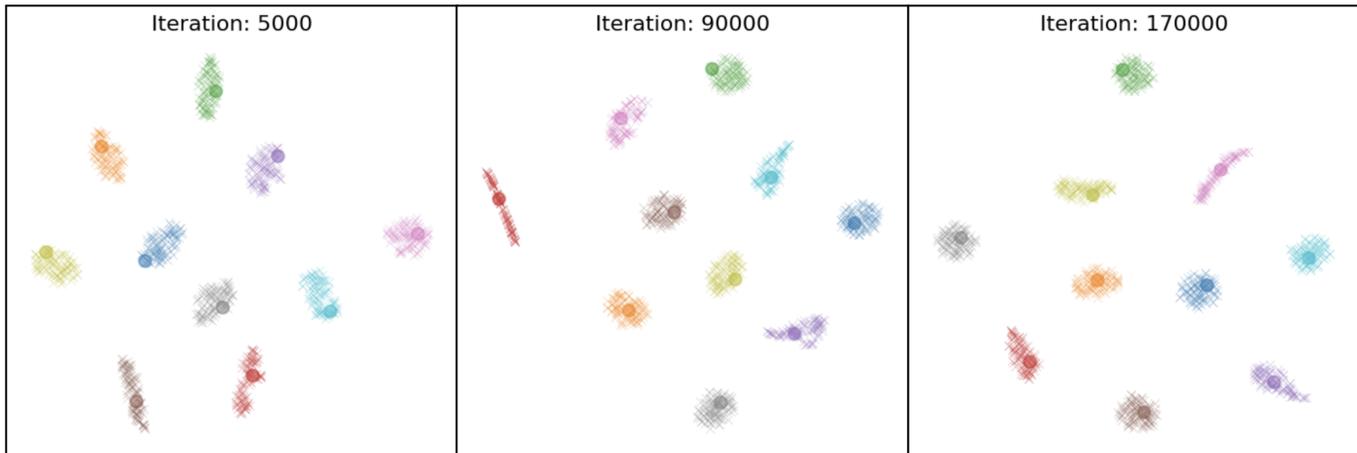


Fig. 4: We randomly sample a data point from each category in the test set of CIFAR10. Then, 100 different adversarial samples corresponding to the natural images are generated. the latent representations of the natural and their adversarial counterparts are plotted using t-SNE. In the early stages of the robust training, the domain gap between the natural and adversarial samples is visible in the latent space. Later in the robust training, the adversarial samples start grouping around the natural sample such that it is harder for the generator to find adversarial directions.

introduce adversarial robustness to the object recognition task without the degradation in the natural accuracy.

#### ACKNOWLEDGMENT

This work is supported by Bogazici University Research Fund under the Grant Number 17004.

#### REFERENCES

- [1] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, J. Santamaría, M. A. Fadhel, M. Al-Amidie, and L. Farhan, “Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions,” *Journal of big Data*, vol. 8, no. 1, pp. 1–74, 2021.
- [2] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- [3] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, “Practical black-box attacks against machine learning,” in *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*. Association for Computing Machinery, 2017, p. 506–519. [Online]. Available: <https://doi.org/10.1145/3052973.3053009>
- [4] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.
- [5] R. Wiyatno and A. Xu, “Maximal jacobian-based saliency map attack,” *arXiv preprint arXiv:1808.07945*, 2018.
- [6] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, “Deepfool: a simple and accurate method to fool deep neural networks. corr abs/1511.04599 (2015),” *arXiv preprint arXiv:1511.04599*, 2015.
- [7] D. Deb, J. Zhang, and A. K. Jain, “Advfaces: Adversarial face synthesis,” in *2020 IEEE International Joint Conference on Biometrics (IJCB)*, 2020, pp. 1–10.
- [8] K. Ren, T. Zheng, Z. Qin, and X. Liu, “Adversarial attacks and defenses in deep learning,” *Engineering*, vol. 6, no. 3, pp. 346–360, 2020.
- [9] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” *arXiv preprint arXiv:1706.06083*, 2017.
- [10] C. Sitawarin, A. N. Bhagoji, A. Mosenia, M. Chiang, and P. Mittal, “Darts: Deceiving autonomous cars with toxic signs,” *arXiv preprint arXiv:1802.06430*, 2018.
- [11] G. Goswami, N. Ratha, A. Agarwal, R. Singh, and M. Vatsa, “Unravelling robustness of deep learning based face recognition against adversarial attacks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [12] S. G. Finlayson, J. D. Bowers, J. Ito, J. L. Zittrain, A. L. Beam, and I. S. Kohane, “Adversarial attacks on medical machine learning,” *Science*, vol. 363, no. 6433, pp. 1287–1289, 2019.
- [13] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [14] F. Feng, X. He, J. Tang, and T.-S. Chua, “Graph adversarial training: Dynamically regularizing based on graph structure,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 6, pp. 2493–2504, 2021.
- [15] J. Tang, X. Du, X. He, F. Yuan, Q. Tian, and T.-S. Chua, “Adversarial training towards robust multimedia recommender system,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 5, pp. 855–867, 2020.
- [16] F. Tramer and D. Boneh, “Adversarial training and robustness for multiple perturbations,” in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 5858–5868. [Online]. Available: <http://papers.nips.cc/paper/8821-adversarial-training-and-robustness-for-multiple-perturbations.pdf>
- [17] W. Zhou, X. Hou, Y. Chen, M. Tang, X. Huang, X. Gan, and Y. Yang, “Transferable adversarial perturbations,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 452–467.
- [18] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, “Theoretically principled trade-off between robustness and accuracy,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 7472–7482.
- [19] H. Jiang, Z. Chen, Y. Shi, B. Dai, and T. Zhao, “Learning to defend by learning to attack,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 577–585.
- [20] Y. Jang, T. Zhao, S. Hong, and H. Lee, “Adversarial defense via learning to generate diverse attacks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2740–2749.
- [21] A. Jeddi, M. J. Shafiee, M. Karg, C. Scharfenberger, and A. Wong, “Learn2perturb: an end-to-end feature perturbation learning to improve adversarial robustness,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1241–1250.
- [22] O. K. Yüksel and İ. M. Baytaş, “Adversarial training with orthogonal regularization,” in *2020 28th Signal Processing and Communications Applications Conference (SIU)*. IEEE, 2020, pp. 1–4.
- [23] H. Zhang and J. Wang, “Defense against adversarial attacks using feature scattering-based adversarial training,” *arXiv preprint arXiv:1907.10764*, 2019.
- [24] S. Lee, H. Lee, and S. Yoon, “Adversarial vertex mixup: Toward better adversarially robust generalization,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 272–281.
- [25] H. Zhang and W. Xu, “Adversarial interpolation training: A simple

- approach for improving model robustness,” 2020. [Online]. Available: <https://openreview.net/forum?id=Syejj0NYvr>
- [26] J. Wang and H. Zhang, “Bilateral adversarial training: Towards fast training of more robust models against adversarial attacks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6629–6638.
- [27] H. Lee, S. Han, and J. Lee, “Generative adversarial trainer: Defense to adversarial perturbations with gan,” *arXiv preprint arXiv:1705.03387*, 2017.
- [28] P. Samangouei, M. Kabkab, and R. Chellappa, “Defense-gan: Protecting classifiers against adversarial attacks using generative models,” *arXiv preprint arXiv:1805.06605*, 2018.
- [29] Y. Dong, Z. Deng, T. Pang, J. Zhu, and H. Su, “Adversarial distributional training for robust deep learning,” in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 8270–8283.
- [30] H. Wang and C.-N. Yu, “A direct approach to robust deep learning using adversarial networks,” *arXiv preprint arXiv:1905.09591*, 2019.
- [31] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *arXiv preprint arXiv:1406.2661*, 2014.
- [32] A. Genevay, G. Peyré, and M. Cuturi, “Gan and vae from an optimal transport point of view,” *arXiv preprint arXiv:1706.01807*, 2017.
- [33] T. Salimans, H. Zhang, A. Radford, and D. Metaxas, “Improving gans using optimal transport,” *arXiv preprint arXiv:1803.05573*, 2018.
- [34] A. Krizhevsky, “Learning multiple layers of features from tiny images,” *University of Toronto*, 05 2012.
- [35] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with unsupervised feature learning,” in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [36] C. Xiao, B. Li, J.-Y. Zhu, W. He, M. Liu, and D. Song, “Generating adversarial examples with adversarial networks,” *arXiv preprint arXiv:1801.02610*, 2018.
- [37] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [38] “Cifar10 adversarial examples challenge.” [Online]. Available: [https://github.com/MadryLab/cifar10\\_challenge.git](https://github.com/MadryLab/cifar10_challenge.git)
- [39] “Tensorflowsinkhorn.” [Online]. Available: <https://github.com/jaberkow/TensorFlowSinkhorn.git>
- [40] “Feature scattering adversarial training.” [Online]. Available: <https://github.com/Haichao-Zhang/FeatureScatter.git>



**İnci M. Baytaş** is currently an Assistant Professor in the Department of Computer Engineering at Boğaziçi University. She received her Ph.D. degree from the Department of Computer Science and Engineering at Michigan State University in 2019. Her research interests include machine learning, deep learning, adversarial robustness, temporal analysis, and biomedical informatics. She has served as program committee member for AAAI since 2019, and as reviewer for premier journals, such as IEEE Transactions on Knowledge and Data Engineering,

IEEE Transactions on Pattern Analysis and Machine Intelligence, IEEE Journal of Biomedical and Health Informatics, and IEEE Transactions on Services Computing.



**Debayan Deb** He received his PhD degree in the Department of Computer Science and Engineering at Michigan State University in 2021. His research interests include pattern recognition, computer vision, and machine learning with applications in biometrics. He served as program committee member for CVPR and ICCV since 2020, as as reviewer for premier journals, including IEEE Transactions on Information Forensics & Security and IEEE Transactions on Pattern Analysis and Machine Intelligence.