Robust Proxy: Improving Adversarial Robustness by Robust Proxy Learning

Hong Joo Lee and Yong Man Ro, Senior Member, IEEE,

Abstract—Recently, it has been widely known that deep neural networks are highly vulnerable and easily broken by adversarial attacks. To mitigate the adversarial vulnerability, many defense algorithms have been proposed. Recently, to improve adversarial robustness, many works try to enhance feature representation by imposing more direct supervision on the discriminative feature. However, existing approaches lack an understanding of learning adversarially robust feature representation. In this paper, we propose a novel training framework called Robust Proxy Learning. In the proposed method, the model explicitly learns robust feature representations with robust proxies. To this end, firstly, we demonstrate that we can generate class-representative robust features by adding class-wise robust perturbations. Then, we use the class representative features as robust proxies. With the classwise robust features, the model explicitly learns adversarially robust features through the proposed robust proxy learning framework. Through extensive experiments, we verify that we can manually generate robust features, and our proposed learning framework could increase the robustness of the DNNs.

Index Terms—Robust perturbation, class-wise robust perturbation, robust proxy learning.

I. INTRODUCTION

R ECENTLY, Deep Neural Networks (DNNs) have achieved great performance in various machine learning tasks [1]–[3]. Despite the phenomenal success of DNNs, they are highly vulnerable to adversarial attacks [4]–[9]. By adding small and imperceptible perturbation to input data (*i.e.*, adversarial examples), adversarial perturbations can effectively fool DNNs. Such vulnerability of DNNs could lead to security problems and the loss of the reliability of DNNs.

To mitigate the potential threat of adversarial attacks, a number of defenses have been proposed [6], [10]–[17]. Among the existing defenses, Adversarial Training (AT) has been demonstrated to be the most effective defense strategy [18], [19]. It trains DNNs with adversarial examples by solving minmax optimization problems between the adversarial perturbation and model parameters.

However, some recent works started tackling the limitation of existing adversarial training schemes through the lens of feature representation [20]–[25]. They claim that wellgeneralized feature representation could improve the adversarial robustness. To enhance the feature representation, they apply AT framework to a self-supervised or unsupervised learning scheme such as SimCLR [23], [24]. However, they just applied the existing representation learning framework to AT framework, there is a lack of a deeper understanding of adversarially robust feature representation.

In the lens of adversarial robustness, it has been known that the disagreement between standard and adversarial robustness stems from differently trained features representation [26]-[28]. The vulnerability of DNNs arises from naturally learned non-robust feature components, and they are highly correlated with adversarial prediction. Ilyas et al. [26] demonstrated that features consist of robust and non-robust components. They claimed that the adversarial examples are directly attributed to the presence of non-robust components and these nonrobust components are brittle by adversarial perturbations while useful for prediction in the standard setting. Kim et al. [27] explicitly distilled features into the robust and nonrobust components. Specifically, they disentangled features into the robust and non-robust channels. Then, they showed that the vulnerability mainly stems from non-robust channels rather than robust channels. The aforementioned analyses of the adversarial vulnerability commonly argue that the feature representations that are learned to correctly predict and robustly predict are different. Therefore, to improve adversarial robustness, it is necessary to learn adversarially robust feature representations. Although there are many studies to identify the robust and non-robust features, only a few works have been conducted to exploit these robust and non-robust features.

Based on the aforementioned adversarially robust feature representation view, in this paper, we raise the following intriguing, yet thus far overlooked questions:

"How can we make DNNs learn the adversarially robust feature?"

To address the question, in this paper, we generate class representative robust features for all classes. Then, with the generated features, we train DNNs to explicitly learn adversarially robust features.

To generate the robust feature, firstly, we employ a feature distillation method proposed in [27] and distill the feature into robust and non-robust channels. With the distilled features, we quantify the effect of the non-robust channels on the prediction by the gradient of them. Then, we optimize the input to minimize the magnitude of the gradient of non-robust channels. Specifically, we add a Robust Perturbation (RP, r) to input data and optimize the perturbation to reduce the gradient of non-robust channels. After that, we extend the process of optimizing robust Perturbation (CRP, r^k). The CRP is added

This work was conducted by Center for Applied Research in Artificial Intelligence (CARAI) grant funded by DAPA and ADD (UD230017TD).

H. J. Lee, and Y. M. Ro are with the Image and Video Systems Laboratory, School of Electrical Engineering, Korea Advanced Institute of Science and Technology (KAIST), Daejeon, 34141, South Korea (e-mail: dlghdwn008@kaist.ac.kr; ymro@kaist.ac.kr). *Corresponding author: Yong Man Ro.*

to any input data corresponding to the target class k and makes the input robust against adversarial perturbation. To optimize CRP, we exploit the Empirical Risk Minimization (ERM) optimization which is considered a successful recipe for finding classifiers with small population risk. Specifically, we propose a novel optimization process called Class-wise ERM Optimization (CEO). In the CEO algorithm, we quantify the empirical risk of the target class as the expectation of gradient for the non-robust channels. Then, we optimize CRP to reduce the gradient. Through empirical and theoretical analysis, we show that CRP makes the corresponding class input not easily attacked by adversarial perturbation and provides a robust prediction.

With the optimized CRP, we generate a class-representative robust feature called Robust Proxy. To generate a robust proxy, we randomly sample each class image from the training dataset and add CRP to the corresponding class images. Then, we extract features from the CRP-added images and use the features as robust proxies. During the training, DNNs explicitly learn the representation of robust proxy through the proposed robust proxy learning framework. For each proxy, we pull the data of the same class close to the proxy and push others away in the feature space, allowing the model to explicitly learn adversarially robust features.

The major contributions of the paper are as follows.

- We propose a novel way to generate adversarially robust features by optimizing robust perturbation. Then, we extend the robust perturbation to class-wise robust perturbation to generate class-representative robust features.
- With the CRP, we train DNNs with the proposed learning framework called Robust Proxy Learning. In the proposed method, we train the DNNs to explicitly learn adversarially robust features by using robust proxies.
- Through extensive experiments, we show that we could explicitly learn robust features and improve the robust-ness.

II. RELATED WORK

A. Understanding Adversarially Robust Features

As adversarial vulnerability has attracted significant attention, many works are devoted to getting to the bottom of the vulnerability [26]–[28]. An early study tended to view adversarial examples as a result of the excessive linearity nature of DNNs in high-dimensional spaces [4]. In another study, it has been regarded as statistical fluctuations in the data manifold [7], [29].

Recently, a new perspective on the phenomenon of adversarial vulnerability is proposed [26]–[28], [30]. In contrast to previous studies, these works figure out the vulnerability as a view of feature representations. Tsipras et al. [30] showed that the goals of learning features for standard accuracy and adversarial robustness might be at odds. Specifically, they argue that the features for adversarial robustness and for standard accuracy are fundamentally different. Then, the features learned by robust models tend to align better with salient data characteristics and human perception. Ilyas et al. [26] demonstrated that adversarial vulnerability is a consequence of nonrobust components of features. These non-robust components are useful and highly predictive for standard performance, yet easily broken by adversarial perturbations. In contrast, robust components still can provide robust prediction results even with adversarial perturbation. Kim et al. [27] explicitly distilled features into robust channels and non-robust channels. Then, they showed that non-robust channels are directly related to adversarial predictions. Also, in [27], they addressed that the robust channels are robust on the noise variation and invariant to the existence of the adversarial perturbation. In contrast, the non-robust channels are brittle and easily change the model prediction by noise variation. In the rest of the paper, we use the definition of robust/non-robust features defined in [27].

The aforementioned studies commonly argued that features for a standard performance and for adversarial robustness are different. Also, the vulnerability stems from the non-robust feature components that are brittle and incomprehensible to humans. Therefore, it is necessary to learn robust feature representations. In this context, recently some works try to improve the robustness by exploiting the robust feature representations. Yang et al. [31] proposed the Deep Robust Representation Disentanglement Network (DRRDN) model to disentangle the class-specific representation and class-irrelevant representation. To this end, they employed a disentangler to extract and align the robust representations from both adversarial and natural examples. With the disentangler, they eliminate the effect of adversarial perturbations and improve the robustness. Kim et al. [32] proposed a way to extract robust and non-robust features based on causality. They demystified causal features on adversarial examples in order to uncover inexplicable adversarial origins through a causal perspective. To this end, they proposed adversarial instrumental variable (IV) regression as a means to identify the causal features pertaining to the causal relationship of adversarial prediction on adversarial examples. Then, they improve the robustness by exploiting the causal features.

In this work, we propose a new approach that explicitly learns robust feature representations rather than heuristically employing existing learning algorithms.

B. Adversarial Training

Adversarial Training (AT) is one of the most effective approaches to defending against adversarial attacks [6], [15]– [17], [33], [34]. By solving a min-max optimization between model parameters and adversarial perturbation, it improves the adversarial robustness. Madry et al. [6] proposed a PGD-based adversarial training method and achieved the first empirical adversarial robustness. Since then, it became a milestone in adversarial training methods.

AT by Enhancing Feature Representation: Many recent studies tackled the limitation of existing adversarial training methods in the lens of feature representation [20]–[25]. Then, they started to study how to improve feature representation by exploiting self-supervised / unsupervised learning scheme. They tried to apply AT framework to a self-supervised / unsupervised pretraining task to make DNNs learn robust data representation. Mao et al. [20] empirically analyze the feature representations under adversarial attack and showed

that adversarial perturbations shift the feature representations of adversarial examples away from their true class and closer to the false class. With the empirical observations, they employ triplet-wise distance loss in the AT framework and improve the robustness. Fan et al. [21] proposed a unified adversarial contrastive learning framework that learns well transferable feature representations. Similar works have been done to enhance the feature representation against adversarial attacks [22], [23], [35], [36].

However, since they heuristically employ existing representation learning framework to AT framework, they do not take any consideration of the adversarially robust feature representations into account. It has been known that naturally learned feature representation that aims to correctly predict is different from adversarially robust feature representation. Therefore, it is necessary to consider the adversarially robust feature representation to improve the adversarial robustness.

III. GENERATING CLASS REPRESENTATIVE ROBUST FEATURES

The main contribution of this paper is to generate classrepresentative robust features. Then, we train the model robustly by using the class representative robust features. In this section, we first describe how to generate class-representative robust features. Specifically, to clarify how to distill robust and non-robust channels from features, we briefly revisit [27]. Then, with the distilled features, we explain how to generate adversarially robust features and class representative robust features. Through the proof concept experiments, we verify that we could generate class-representative robust features by using Class-wise Robust Perturbations (CRP).

A. Revisiting Robust & non-Robust Feature Distillation

In this section, we revisit how to distill features into robust and non-robust channels. In this paper, the definition of robust and non-robust channels stems from [27]. Let z indicates the intermediate feature of the model f such that $z = f_l(x)$, where $f_l(\cdot)$ is *l*-th layer output of the given model. Then, $f_{l+}(\cdot)$ represents subsequent network after the *l*-th layer. Therefore, the prediction of the model can be written as $y' = f_{l+}(f_l(x))$. Since the last convolution layer contains higher-level features and more discriminative features than the earlier convolution layer, we identify the robust features at the last convolutional layer before the global average pooling layer or fullyconnected layer. Here, the feature z has C channels and each channel has inherent feature variation. The σ_z indicates inherent feature variation of the feature z for each channel. Therefore, it can be written as $\sigma_z = [\sigma_{z,1}, \sigma_{z,2}, \dots, \sigma_{z,C}]$ (given parameters). With the given parameter (σ_z), we set the criterion for comparison as $T = max(\sigma_z^2)$, where T denotes the maximum tolerance of the noise variation of the original feature. Here, we use T as the threshold that discriminates robust and non-robust channels. Then, we find the noise variation σ to estimate the prediction sensitivity of each feature channel along the noise intervention. If the noise variation $\sigma_c^2 > T$, the channel is regarded as a robust channel. To find the noise variation (σ), in [27], they exploited information

bottleneck [37]. According to the definition of information bottleneck in [27], [37], it could find maximally informative representation for target labels, restraining input information, concurrently. Therefore, by using the information bottleneck, we can quantify the feature importance and information flow for the target labels. Then, we can quantify the importance of each feature by utilizing the information bottleneck. The optimization objective can be written as follow:

$$\min_{\sigma} L = \underbrace{-y \cdot \log(f_{l+}(f_l(x) + \sigma \cdot \epsilon))}_{\text{cross-entropy}} + \underbrace{\beta D_{KL}[p(z|x)|q_{\sigma}(z))]}_{\text{KL divergence}}.$$
(1)

The first term indicates cross-entropy and the second term indicates KL divergence between the original feature and noise added feature. Also, the ϵ indicates Gaussian noise, and y denotes the ground-truth label. Through the optimization process, we find noise variation $\sigma = [\sigma_1, \sigma_2, \ldots, \sigma_C]$ (optimized parameters). In [27], they have analyzed that σ is related to variance in Gaussian, thus large σ can allow for large variation capacity in the channel, which makes it have the ability to overcome feature variation. On the other way, small σ only allows for a small variation capacity in the channel, which makes it brittle to feature variation.

After we optimize σ , following the aforementioned criterion, we find the robust channel index $(i_r = [i_{r,1}, i_{r,2}, \ldots, i_{r,C}])$. If the optimized noise variation $\sigma_c^2 > T$, the channel is regarded as robust channel and $i_{r,c} = 1$. Then, the non-robust channels are simply reversed from the robust channel index such that $i_{nr,c} = 1 - i_{r,c}$. Finally, the set of robust channel features can be written as $z_r = z \cdot i_r$. Similarly, the set of non-robust channel feature z_{nr} can be written as $z_{rr} = z \cdot i_{rr}$. In this way, the z can be expressed as $z = i_r \cdot z + i_{nr} \cdot z$.

B. Generating Robust Features

1) Problem Definition: In this section, we describe how to generate robust features. To generate robust features, we generate robust inputs and use the features extracted from the inputs as robust features. To this end, we manipulate the input by adding perturbation and define the perturbation as Robust Perturbation (RP, r), which makes the input robust against adversarial perturbation. Then, the input is regarded as the robust input if the model prediction maintains its original prediction, even though there exists adversarial perturbation. Therefore, the problem can be defined as follows,

Making
$$f(x + r + \delta^*) \approx f(x)$$
 by optimizing r
where $\delta^* = \underset{||\delta|| < \epsilon}{\operatorname{argmax}} \mathcal{L}(f(x + r + \delta)), y),$ (2)

where δ^* is an adversarial perturbation that attacks x + r, ϵ is adversarial perturbation budget, x is an input image, and \mathcal{L} denotes the objective function such as cross-entropy. The equation can be interpreted that even though the input x + r is attacked by the adversarial perturbation, the prediction is maintained and the feature extracted from x+r can be regarded as a robust feature. In the following section, we will explain how to optimize r in detail.

Algorithm 1: Robust Perturbation Generation Input: Pre-trained classifier $f(\cdot)$, loss function \mathcal{L}_{robust} , input image x Output: Robust control cont		TABLE 1 ROBUSTNESS COMPARISON WHEN ROBUST PERTURBATION r is addee input data. Note that the adversarial perturbation is generated by PGD-20 on $x_i + r_i$.				
for T iterations do \triangleright Add robust perturbation to		Madry	$x_i \\ x_i + r_i$	46.5 75.3	20.2 48.7	
$x' \leftarrow x + r$ input image $r \leftarrow r - \alpha \nabla_r (\mathcal{L}_{robust}(x', r))$ $\text{low of the second of the secon$	TRADES	$x_i + r_i$	48.8 77.91	21.3 46.5		
	decent $(\alpha: lr)$	MART	$x_i + r_i$	49.1 79.8	21.2 49.2	
return r			x_i	52.1	21.6	

2) Optimizing Robust Perturbation: Recall that the adversarial vulnerability mainly stems from non-robust components of feature [26], [27], [30]. Specifically, in [27], the vulnerability is highly correlated with the non-robust channel of z. If the non-robust channels have a large impact on the prediction, it can be interpreted that the input is vulnerable to adversarial perturbation. Therefore, it is necessary to reduce the influence of non-robust channels on prediction. To this end, in this paper, we quantify the influence of non-robust channels by measuring the gradient of the cost function with respect to nonrobust channels and reducing them by optimizing r. Note that, the gradient of the non-robust channels represents how small changes at each channel affect the prediction. The gradient of non-robust features can be described as follows:

$$\mathcal{G}_{nr} = \frac{\partial}{\partial z_{nr}} \mathcal{L}_{base}(f(x+r), y), \tag{3}$$

where z_{nr} is the set of non-robust channels in feature z = $\{z_r \cup z_{nr}\}, f(\cdot)$ is a model prediction, and \mathcal{L}_{base} is a loss function to ensure correct prediction. We define the loss function as $\mathcal{L}_{base} = -c \cdot \max(\max_{i \in \mathcal{L}_{i}} (f(x)_{i}) - f(x)_{i}, 0)$ due to its empirical effectiveness of optimization performance in [6]. c determines the trade-off between the size of the perturbation added to the input and the degree of correct prediction in Eq.5. \mathcal{G}_{nr} is a quantified value that quantifies how much a change in the non-robust feature changes the correct prediction. If \mathcal{G}_{nr} has a large value, it indicates that we can easily change the prediction for the ground-truth class. The gradient of nonrobust channel can be simplified as follows:

$$\mathcal{G}_{nr} = \frac{\partial}{\partial z_{nr}} \mathcal{L}_{base}(f(x+r), y)$$

$$= \frac{\partial z}{\partial z_{nr}} \frac{\partial}{\partial z} \mathcal{L}_{base}(f(x+r), y)$$

$$= i_{nr} \cdot \frac{\partial}{\partial z} \mathcal{L}_{base}(f(x+r), y).$$

$$(\because z = i_r \cdot z + i_{nr} \cdot z)$$
(4)

Following Eq.4, we could simply reduce the gradient of non-robust channels by multiplying non-robust channel index (i_{nr}) to the gradient of objective function respect to the feature $(\frac{\partial}{\partial z} \mathcal{L}_{base}(\cdot))$. Using the gradient, we optimize robust perturbation by optimizing the following objective:

$$\mathcal{L}_{robust} = \mathcal{L}_{base}(f(x+r), y) + \|\mathcal{G}_{nr}\|_2 + \|r\|_2.$$
 (5)

то

Model	Input Types	CIFAR-10	Tiny-ImagNet
Madry	$x_i \\ x_i + r_i$	46.5 75.3	20.2 48.7
TRADES	$x_i + r_i$	48.8 77.91	21.3 46.5
MART	$x_i + r_i$	49.1 79.8	21.2 49.2
HELP	$x_i \\ x_i + r_i$	52.1 78.5	21.6 48.4

Optimizing Eq. 5 means that we find a small and imperceptible perturbation that makes the model predict well and reduces the gradient of non-robust features. Then, reducing the gradients of non-robust channels contains the same effect of reducing σ_{nr} to resist adversarial perturbation. Therefore, features extracted from those inputs (x+r) can be a robust feature. Note that there is no regularization parameter before the gradient norm of non-robust features and the norm of perturbation in Eq.5 The optimization algorithm is described in Algorithm 1.

3) Robustness Analysis with Robust Perturbation: The goal of optimizing r is to make the input itself robust against adversarial perturbation. If we can robustify the input by adding r to the input, we can extract features from these inputs and regard them as robust features. For verification, in this section, we conduct a proof concept experiment to verify whether we can make the input itself robust against adversarial attacks by augmenting the input.

In the proof concept experiment, we optimize r_i for all corresponding images x_i in the test dataset of CIFAR-10 and Tiny-ImageNet on the pre-trained AT models (Madry [6], TRADES [16], MART [15], and HELP [17]). The results are described in Table I. Table I shows the robustness comparison under the PGD-20 attack according to different input types. In the table, x_i denotes an accuracy when the adversarial perturbation is added to the original input image, and $x_i + r_i$ denotes an accuracy when the adversarial perturbation is added to $x_i + r_i$. Note that to verify that r_i truly makes the input robust, we generate the adversarial perturbation on $x_i + r_i$. As shown in the table, we verify that adding r_i to input data can significantly improve the robustness. Furthermore, even though the adversarial perturbation is generated on $x_i + r_i$, r_i could successfully improve the robustness (robust against adaptive attack). This can be interpreted that r_i does not cause gradient obfuscation [9], and reducing the gradient of the nonrobust channels makes the input robust itself.

C. Optimizing Class-wise Robust Perturbation

We have analyzed that augmenting input with robust perturbation makes the input robust against adversarial perturbation. In this section, we expand the concept of r to Class-wise Robust Perturbation (CRP, r^k) where k denotes the class index. The CRP is a class-specific perturbation that can be applied to any corresponding class input (x^k) and improve



Fig. 1. Overview of the proposed CEO process. The CRP is added to input for the corresponding class, fed into the freezed model, and extract feature. Then, the features are distilled into robust / non-robust channels and we measure the empirical risk of non-robust channels to optimize the CRP. The calculated objective function is backpropagated to CRP and we only optimize CRP.

the robustness. For example, if we generate a CRP for the dog class (r^{dog}), it can be applied to any dog class images and improve the robustness. To this end, we propose a novel optimization method called Class-wise ERM Optimization (CEO) by extending Empirical Risk Minimization (ERM) algorithm, which is considered a successful recipe for finding classifiers with small population risk [6]. In order to make r^k have universality in corresponding class inputs, we measure the empirical risk of non-robust channels for corresponding class images. Let $X = \{x_1^k, x_2^k, ..., x_N^k\}$ be a subset of class-k images sampled from the training data and the feature extracted from each input is z_n^k . Then, the empirical risk of non-robust channels can be formulated as follows:

$$\mathcal{L}_{feature} = \frac{1}{N} \sum_{n=1}^{N} \left\| \frac{\partial}{\partial z_{n,nr}^{k}} \mathcal{L}_{base}(x_{n}^{k} + r^{k}, y_{n}^{k}) \right\|, \quad (6)$$

where N denotes the number of images in the target class and $z_{n,nr}^k$ denotes the set of non-robust channels of n-th image. Therefore, the total objective function for CEO is

$$\mathcal{L}_{CEO} = \frac{1}{N} \sum_{n=1}^{N} \mathcal{L}_{base}(x_n^k + r^k, y_n^k) + \mathcal{L}_{feature}.$$
 (7)

The objective function means that the CRP is optimized to correctly classify the target class k and reduce the effect of the non-robust channels within the class. Here, the gradient of CEO with non-robust channels is gradually converged to local optima, where it asymptotically closes to zero vector to make it be small population risk. Through the CEO, we could generate CRP that could improve the robustness of the target class and generate class representative common robust features. Fig. 1 gives an overview of how to optimize the CRP.

1) Analysis of CRP: The goal of optimizing r^k is to generate the class-wise robust perturbation that could improve the robustness of the corresponding target class images. For verification, we also conduct proof concept experiments. In the experiments, we optimize r^k for each class from the training set and apply them to the test set images. Table II shows the robustness comparison under the PGD-20 attack according to different input types. In the table, x_i^k denotes an accuracy

TABLE IIROBUSTNESS COMPARISON WHEN CLASS-WISE ROBUST PERTURBATION r^k is added to input data. Note that the adversarialPERTURBATION IS GENERATED BY PGD-20 on $x_i^k + r^k$.

Model	Input Types	CIFAR-10	Tiny-ImagNet
AT	$x_i^k \overset{x_i^k}{+} r^k$	46.5 69.2	20.2 35.3
TRADES	$x_i^k + r^k$	48.8 69.8	21.3 35.2
MART	$x_i^k + r^k$	49.1 68.27	21.2 37.0
HELP	$x_i^k x_i^k + r^k$	52.1 70.22	21.6 36.8



Fig. 2. Visualization results of CRPs (r^k) and CRPs added images $(x^k + r^k)$. k denotes car and dog classes.

when we use the original input image, and $x_i^k + r^k$ denotes an accuracy when adding CRP to corresponding target class images. Different from the result in Table I that generates robust perturbation for all corresponding input images, in this experiment, we generate one CRP per class and applied it to all images that correspond class. The results are shown in Table II. As shown in the table, we verify that adding CRP to input also significantly improves the robustness. In other words, once we optimize CRP from training images, we can apply CRP to any test image that corresponds to the target class and robustify the target class images. Furthermore, we



Fig. 3. Distribution of L2-norm of non-robust features. (a) and (b) is the non-robust feature gradient of CIFAR-10 and Tiny-ImageNet respectively.



Fig. 4. Feature similarity distribution between positive samples. The blue bar denotes the feature similarity distribution between images belonging to the same class. The orange bars denote the feature similarity distribution between the same classes when CRP is added to the image.

visualize the CRPs.Fig. 2 shows the examples of CRPs (r^k) and CRPs added images $(x^k + r^k)$. As shown in the figure, CRP contains the semantic information of the corresponding class, and its magnitude is very small. Therefore, even if we add the CRPs to corresponding class images, we can robustify the corresponding image well without any significant change.

Measuring the Gradient of Non-robust Channels: As we discussed above, the norm of the gradients with respect to the non-robust channels are related to adversarial vulnerability. Then, we have shown that adding CRP to input can improve the adversarial robustness. In this section, we verify whether the increase in robustness is caused by the decrease in the gradient of the non-robust channels. To this end, we apply CRPs to input data and measure the L2-norm magnitude of non-robust channels. The results are shown in Fig. 3. As shown in the figure, the gradients of non-robust channels are reduced significantly compared with 'without CRP'. Therefore, we demonstrate that CRPs could improve the robustness by reducing the gradient of non-robust features.

Feature Similarity Measurement: Furthermore, to verify whether adding CRP can generate a feature that can represent the class or not, we measure the feature similarity between samples belonging to the same class (positive samples). The result is shown in Fig. 4. As shown in the figure, when applying the CRP to input data (orange color bar) the similarity between positive samples is increased and has narrow distribution. The experimental result shows that if the features are extracted by adding CRP to any input corresponding to the target class, the features have similar representations and can be regarded as robust and class-representative features.



Fig. 5. Graphical explanation of the proposed robust proxy learning framework. Each robust proxy is colored in black and different shapes denotes different classes. Note that for simplicity, in the case of proxy-negative pair, we only describe for one proxy.

Algorithm 2: Robust Proxy Learning
Input: Classifier $f(\cdot)$, adversarial loss \mathcal{L}_{AT} ,
robust proxy loss \mathcal{L}_{proxy} , extracted feature $z(x)$,
set of all proxies P , training dataset \mathcal{D} ,
Interval T, learning rate τ
Output: Model parameter θ
Initialize θ
for epoch do
if epoch%T=0 do
for $k = 1, 2, 3, \dots, K$ do
Generate class-wise robust perturbation r^k
Randomly sample x^k for each class k , $x^k \subset D$
$p^k \leftarrow z(x^k + r^k)$
end for
for minibatch $X \subset \mathcal{D}$ do
$\mathcal{G}_{\theta} \leftarrow \mathbb{E}_{(x,y)\in B}[\nabla_{\theta}\mathcal{L}_{total}(x,y,p)]$
$\theta \leftarrow \theta - \tau \mathcal{G}_{\theta}$
end for
end for
return θ

IV. PROPOSED ROBUST PROXY LEARNING

A. Generating Robust Proxy

In this section, we introduce our novel robust proxy learning framework. The main idea is to regard each robust proxy as an anchor in Triplet loss and associate it with entire data in a batch. To this end, we first explain how to generate robust proxies. In the previous section, we verified that even if features are extracted by adding r^k to any input belonging to the target class k, the features have similar values. Therefore, we randomly select the input corresponding to each class, add the corresponding CRP, and then extract the features. The extracted features can be regarded as class-representative robust features and we use them as robust proxies.

B. Training with Robust Proxies

With the robust proxies, we make the DNNs explicitly learn adversarially robust features. Fig. 5 briefly illustrates how to learn robust features during the training. In the figure, each proxy is colored in black, and different shapes indicate different classes. As shown in the figure, positive samples and

TABLE III Adversarial robustness comparison on CIFAR-10 and Tiny-ImageNet dataset under white-box attack setting. The backbone models are ResNet-18 and Wide-ResNet28-10 (WRN28-10).

Backbone	Method			CIFAR-10			Tiny-ImageNet				
Buencone	memou	Clean	FGSM	PGD	CW	AA	Clean	FGSM	PGD	CW	AA
	Standard	92.1±0.10	18.1±0.13	$1.1{\pm}0.02$	$0.0{\pm}0.00$	$0.0{\pm}0.00$	$60.3 {\pm} 0.02$	$6.2 {\pm} 0.10$	$0.8{\pm}0.01$	$0.0{\pm}0.00$	$0.0{\pm}0.00$
ResNet-18	Madry Madry+Proxy	83.8±0.02 84.5±0.02	60.3±0.02 62.0±0.04	50.1±0.04 53.1±0.06	48.3±0.03 49.4±0.04	47.0±0.03 48.9±0.04	46.1±0.02 46.3±0.04	23.3±0.05 24.1±0.02	20.2±0.06 22.7±0.03	18.2±0.04 20.5±0.03	17.6±0.05 19.9±0.05
	TRADES TRADES+Proxy	83.0±0.02 84.0±0.04	61.4±0.05 63.1±0.05	53.0±0.06 55.3±0.04	48.6±0.05 49.7±0.05	47.0±0.05 49.8±0.05	47.8±0.06 47.5±0.07	24.4±0.07 25.1±0.05	21.3±0.02 22.5±0.06	18.4±0.04 21.2±0.06	18.2±0.06 20.0±0.05
	MART MART+Proxy	82.3±0.06 84.0±0.05	60.6±0.05 63.4±0.04	54.1±0.08 56.3±0.06	48.5±0.06 51.0±0.05	47.2±0.06 49.9±0.05	47.6±0.07 47.8±0.04	23.9±0.02 25.5±0.06	21.2±0.04 23.1±0.05	19.2±0.03 21.3±0.04	18.1±0.03 20.0±0.02
	HELP HELP+Proxy	84.0±0.04 85.0±0.06	61.9±0.03 63.8±0.06	52.0±0.07 55.3±0.04	49.8±0.05 52.0±0.04	48.6±0.06 51.2±0.04	48.0±0.07 49.0±0.08	24.0±0.05 25.3±0.02	21.6±0.03 23.0±0.05	19.3±0.03 20.5±0.01	17.7±0.05 19.1±0.01
	Standard	$96.2{\pm}0.12$	$22.3{\pm}0.10$	$3.5{\pm}0.01$	$0.0{\pm}0.00$	$0.0{\pm}0.00$	$63.1{\pm}0.02$	$9.2{\pm}0.02$	$2.7{\pm}0.00$	$0.0{\pm}0.00$	$0.0{\pm}0.00$
WRN28-10	Madry Madry+Proxy	85.5±0.02 86.8±0.02	62.3±0.04 64.5±0.05	54.2±0.05 57.8±0.05	50.8±0.06 52.3±0.06	49.9±0.04 51.7±0.06	48.6±0.03 50.0±0.06	25.1±0.08 27.2±0.04	23.0±0.06 25.9±0.05	20.0±0.06 22.3±0.05	18.7±0.06 20.1±0.04
	TRADES TRADES+Proxy	85.2±0.04 86.8±0.05	63.4±0.04 66.2±0.06	56.2±0.06 58.3±0.03	50.7±0.04 53.8±0.02	49.8±0.05 52.0±0.02	50.6±0.07 52.0±0.03	26.8±0.08 28.3±0.02	25.1±0.02 27.7±0.03	21.9±0.03 22.8±0.02	19.0±0.01 21.3±0.03
	MART MART+Proxy	85.7±0.09 86.7±0.08	63.5±0.05 65.6±0.06	56.2±0.06 59.4±0.04	52.4±0.05 54.0±0.04	51.0±0.05 53.1±0.03	50.4±0.04 52.5±0.05	28.2±0.04 30.8±0.05	26.2±0.06 28.7±0.04	23.4±0.03 24.0±0.02	20.4±0.03 22.6±0.02
	HELP HELP+Proxy	86.8±0.08 86.0±0.05	65.0±0.06 67.2±0.04	56.8±0.06 57.7±0.04	53.9±0.05 55.3±0.04	52.5±0.06 54.0±0.03	51.6±0.05 53.1±0.06	27.8±0.04 29.6±0.05	26.7±0.06 27.9±0.05	23.5±0.03 24.2±0.04	20.0±0.03 22.8±0.04

corresponding positive proxy (p^+) are trained to reduce the distance from each other, whereas negative samples and each proxy are trained to increase the distance from each other. This can be formulated as follows:

$$\mathcal{L}_{proxy} = \underbrace{\frac{1}{|P^+|} \sum_{p \in P^+} \sum_{z \in z_p^+} (d(z, p) - m)}_{\text{Pulling to the Proxy}} - \underbrace{\frac{1}{|P|} \sum_{p \in P} \sum_{z \in z_p^-} (d(z, p) + m)}_{\text{Pushing from the Proxy}}$$
(8)

where m denotes a margin, P denotes the set of all proxies, P^+ denotes the set of positive proxies of data in a batch, and $d(\cdot)$ denotes a distance function. Also, z_p^+ denotes the set of positive features of corresponding $p,\,z_p^-$ denotes the set of negative features of p^+ , and we use cosine similarity between two vectors as a distance function. To boost the robustness of existing AT methods, we jointly optimize the existing AT losses and proxy loss. Therefore, we optimize $\mathcal{L}_{total} = \mathcal{L}_{AT} + \mathcal{L}_{proxy}$ to train the model. \mathcal{L}_{AT} denotes the existing loss function of AT methods. For example, in the case of MART [15] the \mathcal{L}_{AT} can be written as \mathcal{L}_{AT} = $BCE(f(x+\delta), y) + \lambda(1-f_y(x))KL(f(x+\delta), f(x)),$ where $f_u(x)$ denotes the output probability of ground-truth class, λ is a tunable scaling parameter that balances the two parts of the loss, and $BCE(\cdot)$ adds the cross-entropy loss and margin loss terms to improve the decision margin of the classifier. Then, \mathcal{L}_{proxy} makes the model learn features in which the effects of non-robust features are suppressed during training and separate the different classes. After training, we do not need robust perturbation for inference. Therefore, during the inference, robust perturbations are not added to the test data. **Training Details:** In this section, we describe the training details of the proposed robust proxy learning. Note that, since robust/non-robust features can be distilled from the adversarially trained model, we initialize the model parameters

of the AT model. Then, we fine-tuned the model parameter with \mathcal{L}_{total} and set the margin value as 1.0. Since generating robust proxies for every epoch increase the training time, we refresh the proxies for every T = 5 epoch. The details algorithm is described in Algorithm 2.

V. EXPERIMENTS

A. Experiment Setting

1) Dataset and Network: We conduct experiments to verify the effectiveness of the proposed method on two datasets (CIFAR-10 [38] and Tiny-ImageNet [39] datasets). The CIFAR-10 dataset consists of 50,000 training images and 10,000 test images with 10 classes. The Tiny-ImageNet dataset has 200 classes and each class has 500 training images, 50 validation images, and 50 test images. For both datasets, we use the ResNet-18 [40] and WideResNet28-10 [41] networks.

2) Attack Settings: To evaluate the defensive performance of the proposed method, we conduct four adversarial attack methods used as benchmarks for evaluating adversarial robustness (FGSM [4], PGD [6], CW [9], and AutoAttack [42]). For all attack methods, we set the perturbation budget $\epsilon = 8/255$. We generate adversarial perturbation with 20 iterations with the step size $\epsilon/10$ for the PGD attack. In the case of CW [9] attack, we use L_{∞} -norm bounded attacks with 200 iterations. In the case of Auto Attack (AA) [42], it includes four attack methods (APGD-CE, APGD-DLR [42], FAB [43], and Square Attack [44]) and generates adversarial perturbation by ensembling them. For FAB attack hyper-parameters, we optimize the perturbation with 100 iterations and 5 random restarts. In the case of Square attack, we fed 5000 queries for the black-box attack.

B. Robustness Evaluation

1) White-box Evaluation: To evaluate the effectiveness of the proposed method, we apply our proposed robust proxy

learning to existing AT methods. Table III shows the robustness evaluation results under various white-box attack settings on two datasets. In the table, 'Clean' denotes the results when testing with a clean dataset (test without perturbation), and 'Standard' denotes the results when testing on the clean model. Note that the clean dataset denotes the original test image that any perturbations such as adversarial perturbations or class-wise robust perturbations are not added. Furthermore, the clean model denotes the model trained only with the clean dataset. In other words, the clean model is trained without adversarial examples and only uses original training data for training. As shown in the table, the clean model (Standard) shows better performance on the clean dataset compared with the adversarially trained model. However, in the case of adversarial robustness (FGSM, PGD, CW, and AA), it cannot defend against adversarial attacks.

In terms of the adversarially trained model, the baseline results are Madry [6], MART [15]¹, TRADES [16]², and HELP $[17]^3$. Then, the results of our proposed proxy learning are AT+Proxy, TRADES+Proxy, MART+Proxy, and HAT+Proxy. We report the mean and variance from 5 checkpoints. As shown in the table, in the case of Madry, when trained with the proposed Robust Proxy Learning, the robustness is improved to 3.5, 3.1, 1.8, and 2.9 in FGSM, PGD, CW, and AutoAttack, respectively, when trained with the CIFAR-10 dataset on ResNet-18. Also, when we combined proxy learning with other AT methods, we can get similar results. It can be interpreted that the proposed robust proxy learning framework can be well adapted to existing AT frameworks and could improve the robustness. Furthermore, the reason why the proposed method can improve robustness is that it explicitly learns robust feature presentation while learning discriminative features between classes through robust proxy learning.

In terms of clean accuracy, the proposed method does not sacrifice clean accuracy. According to [26], the robust features are useful and highly related to the target class. Therefore, exploiting the robust features does not hurt clean accuracy. Furthermore, through the proposed proxy loss (Eq. 8), the model learns class-discriminative representations. With the proxy loss, we pull the data of the same class close to the proxy and push others away in the feature space, allowing the model to explicitly learn class discriminative features. Therefore, the proposed method shows better performance even on clean samples.

2) Black-box Evaluation: Different from white-box adversarial attacks, black-box attacks generate adversarial perturbation from an unknown model. To show the robustness of the proposed proxy learning under the black-box attack settings, we pre-train WideResNet-34-10, then generate adversarial examples from the model. Then, we use the examples to our WideResNet-28-10 models. The black-box attack results are shown in Table IV and V. 'Base' denotes the reimplementation results with existing AT methods, and 'Ours' denotes the results of proxy learning. Table IV shows the black-box result

Metho	od	FGSM	PGD-20	CW	AutoAttack
Madry	Base	81.37	82.41	83.02	82.41
	Ours	82.12	83.89	84.21	83.91
TRADES	Base	82.29	83.01	83.16	82.98
	Ours	83.21	84.39	84.84	84.72
MART	Base	81.76	82.56	83.09	82.59
	Ours	83.37	84.00	84.27	84.21
HELP	Base	82.33	83.03	83.62	83.81
	Ours	84.16	85.60	85.42	85.88

TABLE V BLACK-BOX ATTACK EVALUATION ON TINY-IMAGENET DATASET. THE PERTURBATION IS GENERATED ON WIDEREsNet-34-10.

Metho	d	FGSM	PGD-20	CW	AutoAttack
Madry	Base	45.42	46.09	47.36	46.27
	Ours	47.45	47.99	48.35	48.23
TRADESS	Base	46.30	47.23	48.15	48.52
	Ours	48.45	49.03	49.23	49.98
MADT	Base	47.40	48.14	48.27	48.4
MARI	Ours	48.30	49.15	49.61	49.50
HELP	Base	48.91	49.47	49.88	50.01
	Ours	49.71	49.93	50.64	50.98

on the CIFAR-10 dataset. As shown in the table, our method could improve the adversarial robustness of existing methods. Also, compared with the white-box results, we achieve better robustness and show close to the natural accuracy. Similar results are described in Table V conducted on the Tiny-ImageNet dataset.

3) Sanity Check about Gradient Obfuscation: Some works provide a false sense of security by vanishing the gradient. These methods are not considered to provide actual robustness. This phenomenon is called gradient obfuscation (gradient masking) [19]. It has been widely known that gradient obfuscation-based defense strategy does not provide adversarial robustness [19], [45]–[48]. The adversarial robustness must guarantee their robustness under the worst case. However, the gradient obfuscation-based method supposes that the gradient of the model is not exposed and is unknown. Therefore, if the gradient of the model is exposed to the adversary, it cannot guarantee adversarial robustness. To handle this, it is important to verify whether the defense strategy is based on gradient obfuscation or not. To verify that the defense strategy is not a gradient obfuscation, in [19], [45], they provide several sanity checks as follows:

- It should show better robustness against 1-step attacks (*e.g.*, FGSM) than iterative attacks (*e.g.*, PGD).
- It should show better robustness against black-box attacks than white-box attacks.

To verify whether the defense strategy is based on gradient obfuscation or not, many studies have utilized the above points [46]–[49]. Then, many defense strategies failed to pass the sanity check. However, from the above experiments, we verify that our proposed method does not suffer from gradient obfuscation. 1) In Table III, the proposed defense method shows better robustness against the FGSM attack than the PGD

¹https://github.com/YisenWang/MART

²https://github.com/yaodongyu/TRADES

³https://github.com/imrahulr/hat

TABLE VI Robustness comparison results with recently proposed representation learning-based methods on ResNet-18 model with CIFAR-10.

Method	FGSM	PGD	AutoAttack
TLA	60.1	50.5	47.0
RoCL	60.5	51.6	48.8
AdvCL	61.6	53.2	49.8
AGKD-BML	61.3	52.3	50.0
Ours	62.5	55.3	51.2

TABLE VII Comparison results with recently proposed methods that exploit the robust and non-robust features. The experiment was conducted on CIFAR-10 dataset with ResNet-18 model.

Method	FGSM	PGD	CW	AutoAttack
CAFE	60.5	54.5	50.4	48.5
DRRDN	62.1	52.1	51.3	47.9
Ours	63.8	55.3	52.0	51.2

attack. 2) By comparing the result of Table III, Table IV, and Table V, we verified that it shows better robustness under the black box setting than in the white-box setting. Through the sanity check, we demonstrate that the proposed method does not rely on gradient obfuscation.

C. Comparison with Recently Proposed Methods

Many previous works try to enhance feature representation to improve adversarial robustness. In this section, we compare the proposed robust proxy learning method with recently proposed methods that enhance the feature representations (TLA [20], RoCL [23], AdvCL [21], AGKD-BML [36]). In the case of TLA and AGKD-BML, they are fully supervised learning methods that apply existing metric learning frameworks to AT framework. In the case of RoCL and AdvCL, they apply AT framework to a self-supervised or unsupervised learning scheme for pretraining. Then, finetuned in a supervised manner. Table VI shows the comparison results. As shown in the table, our proposed method achieves better robustness than other methods. Since our proposed method explicitly learns adversarially robust feature representation, it shows better robustness.

Furthermore, recently some works try to improve the robustness by exploiting the robust features [31], [32]. Table VII shows the comparison results. DRRDN [31] denotes a method that disentangles the features into class-specific features and class-irrelevant features. Then, exploit class-specific features as robust features. In the case of CAFE [32], it extracts robust features by adversarial instrumental variable (IV) regression. The experiment was conducted on the CIFAR-10 dataset with the ResNet-18 model. As shown in the table, our proposed method shows better robustness than recently proposed methods that exploit robust features.

1) Fine-tune with Proxy Loss: As we referred above, many existing feature representation learning methods exploit their methods as pretraining in an unsupervised/semi-supervised manner. Then, finally, they fine-tune the model in a supervised

TABLE VIII IMPROVING THE ROBUSTNESS BY COMBINING THE PROPOSED METHOD WITH EXISTING REPRESENTATION LEARNING-BASED METHODS

Pretraining	Fine-tune	FGSM	PGD- L_{∞}	AutoAttack	$PGD-L_2 \\ (\epsilon = 0.25)$	$\begin{array}{l} \text{PGD-}L_2\\ (\epsilon = 0.5) \end{array}$
RoCL	Base	60.5	51.6	48.8	70.3	60.8
	Ours	61.5	53.1	50.0	72.5	62.9
AdvCL	Base	61.6	53.2	49.8	71.1	60.0
	Ours	62.8	54.8	51.1	73.2	63.0

manner. Therefore, the proposed method can be combined with existing representation learning as a fine-tuning process. To this end, we pretrain the model with AdvCL [21] and RoCL [23]. Then, we fine-tuned the model with the proposed proxy loss. The results are shown in Table VIII. As shown in the table, our method can be combined with AdvCL and RoCL and improve the robustness over them. Furthermore, we verify the effectiveness of the proposed method under unseen adversarial attack (L_2 -norm bounded attack). As shown in Table VIII, our proposed method still shows better robustness under unseen attack. The results can be interpreted that the proposed method can ensure robustness against unseen types of perturbations. Since we explicitly learned a feature that is resistant to noise variation, it shows robustness against unseen perturbations.

2) Computation Cost for Data Sampling in Proxy Learning: Another advantage of Robust Proxy Learning compared to existing representation learning-based AT methods is that the complexity of training computation is low. Here, the training complexity represents the amount of computation required to solve the entire training dataset. Let N denote the number of samples in the training dataset. Since most of existing methods take a tuple of data as a unit input (anchor, positive, negative), it requires high training complexity. For example, in the case of AdvCL and RoCL that exploit contrastive loss, they take a pair of data as input thus they require $O(N^2)$ training complexity. Furthermore, in the case of TLA that exploits Triplet loss, it takes triplets of data thus it requires $O(N^3)$ training complexity. Compared with these methods, in our method, we generate proxies for each class and compare every proxy with all samples. Therefore, the training complexity of the proposed method is O(NC) where C denotes the number of classes. Since $C \ll N$, the training complexity of the proposed method is much less than others.

D. Attacking by Maximizing Gradient of Non-robust Features

One way to directly attack the proposed method is to maximize the gradient of non-robust channels. To generate adversarial perturbation, we modify the Eq. 5 as follows:

$$\mathcal{L}_{attack} = -\mathcal{L}_{base}(f(x+p), y) - \|\mathcal{G}_{nr}\|_2 + \|p\|_2, \quad (9)$$

where p is an adversarial perturbation that maximizes the gradient of non-robust channels and leads to misclassification. The robustness evaluation result is described in Table IX. In the experiment, we set the perturbation budget as $\epsilon = 0.03$. As shown in the table, when maximizing the gradient of non-robust channels, the baseline models (Madry, TRADES, MART, and HELP) are significantly broken. However, when we train the model with the proposed proxy



Fig. 6. Visualization of learned representation. The 'Intermediate' denotes the visualization results of the intermediate feature, and the 'Robust Channels' denotes the visualization results of a set of channels with rarely be manipulated by adversarial perturbations. The 'Non-robust Channels' denotes the visualization results of the set of channels that can potentially be manipulated by adversarial perturbation.

 TABLE IX

 ROBUSTNESS EVALUATION WHEN THE ADVERSARIAL PERTURBATION IS

 GENERATED BY MAXIMIZING THE GRADIENT OF NON-ROBUST CHANNELS

Dataset	Method	Accuracy
	Madry	20.1
	Madry+Proxy	23.2
	TRADES	24.8
CIEAD 10	TRADES+Proxy	27.1
CIFAR-10	MART	22.3
	MART+Proxy	25.7
	HELP	20.6
	HELP+Proxy	23.8
	Madry	9.2
	Madry+Proxy	12.5
	TRADES	11.9
Tiny ImageNet	TRADES+Proxy	14.1
Tilly-Illagervet	MART	12.7
	MART+Proxy	14.1
	HELP	11.2
	HELP+Proxy	14.3

learning (Madry+Proxy, TRADES+Proxy, etc), it shows better robustness. Since the feature that reduces the effect of the nonrobust channels is explicitly learned, it is not easily attacked even if we maximize them, and it shows better robustness.

E. Effectiveness of Robust Proxy

In this section, we verify the effectiveness of robust proxy by ablation study. To this end, we conduct an ablation study by using different types of proxies in Eq.8. Table X shows the experimental results with the ResNet18 model on CIFAR-10 and Tiny-ImageNet datasets when use different types of proxies. In the table, Madry+Proxy_{normal} denotes the results when using a normal instance as the anchor, Madry+Proxy_{avg} denotes the experimental results when using proxy by averaging all robust features. Also, Madry+Proxy_{ours} denotes the results when using our proposed robust proxies as the anchor. As shown in the table, when using the robust proxy, it shows better robustness than when using normal instances as the anchor. Since the normal instances do not guarantee robust representation, they cannot improve the adversarial robustness.

TABLE X Ablation study by using different types of proxies. The experiment was conducted on the ResNet18 model with CIFAR-10 and Tiny-ImageNet.

Dataset	Method	FGSM	PGD	CW	AA
	Madry + Proxy _{normal}	60.1	50.5	48.1	47.0
CIFAR-10	Madry + Proxy	61.0	51.7	48.5	47.1
	Madry + $Proxy_{ours}$	62.0	53.1	49.4	48.9
	Madry + Proxy _{normal}	23.5	21.1	18.0	17.0
Tiny-ImageNet	Madry + Proxy	23.5	20.2	19.0	18.0
	Madry + Proxy _{ours}	24.1	22.7	20.5	19.9

However, in the proposed method, the robust proxies have robust feature representation, we could improve the robustness by robust proxy learning. Furthermore, when using the average of robust proxies, it shows less robustness compared to using the proposed robust proxy. Since simply averaging the features could lead to different feature representations, it could lead the robustness decrement.

F. Visual Interpretation of Learned Representations

1) Feature Visualization: In this section, we show the advantage of the proposed method from the lens of learned feature representation. To this end, we visualize the features that have high similarity and low similarity to proxies. We visualize the features by [27]. In [27], the visualization results are optimized from random image. It optimizes the random input image so that the features extracted from the optimized input image are similar to the robust or non-robust features. Then, the optimized image becomes the visualization result. Fig. 6 shows the visualization results of robust and non-robust channels for each input. In the figure, following the definition of [27], the 'Intermediate' denotes the visualization results of the intermediate feature, and the 'Robust Channels' denotes the visualization results of a set of channels with rarely be manipulated by adversarial perturbations. 'Non-robust Channels' denotes the visualization results of the set of channels that can potentially be manipulated by adversarial perturbation. As



Fig. 7. t-SNE visualization of representations according to the differently trained model. (a) denotes the t-SNE feature visualization results of the existing method (Madry). (b) denotes the t-SNE feature visualization results of the proposed method (Madry+Proxy).



Fig. 8. t-SNE visualization of representations according to different input types. Adding r or r^k gives a much clearer separation among classes.

shown in the figure, the features of 'High Similarity Samples' have semantic information by themselves, and the features of 'Low Similarity Samples' do not have semantic information by themselves. Specifically, in the case of non-robust channels of high similarity samples, it remains the semantical information of ground-truth classes. Therefore, the prediction does not change even if adversarial perturbations are added. In contrast to, in the case of the negative units of low similarity samples, it contains different semantical information from the ground-truth class. Therefore, it is easily manipulated by adversarial perturbation.

2) Verify the effectiveness of Proposed Proxy Learning by *t-SNE visualization:* In the proposed method, the performance improvement mainly stems from the separation of robust and non-robust features. Different from previous works, since our proposed method explicitly learns class discriminative robust features, it shows better robustness. To verify this, we visualize the features using t-SNE on the CIFAR-10. Fig. 7 shows the t-SNE results. Fig. 7 (a) shows features extracted from the base model (Madry). Fig. 7 (b) shows features extracted from the figure, Fig. 7 (b) shows more discriminative feature distribution and clearer class boundary than (a). This can be interpreted that the proposed method can learn more robust and class discriminative features.

3) Verify the effectiveness of CRP by t-SNE visualization: To further demonstrate the efficacy of CRP, we visualize the features according to the input types using t-SNE on the CIFAR-10 dataset. Fig. 8 shows the t-SNE results. As shown in the figure, when we add r or r^k to input, it shows a



Fig. 9. The accuracy comparison of robust and non-robust features by controlling information flows of features. The blue line denotes the accuracy when predicting only with robust features and the red line denotes the accuracy when predicting only with non-robust features. Note that β regulates the total amount of the information that flows into the features.

much clearer class boundary than using the original images. This can be interpreted that r and r^k make the adversary difficult to successfully attack an image, leading to more robust prediction.

G. Effect of Information Bottleneck in Feature Distillation

In this section, we conduct ablation studies to verify the role of the information bottleneck. To this end, we change the β values in Eq. 1 to analyze how the information bottleneck controls the information flow of the robust and non-robust features. In the equation, β regulates the total amount of the information that flows into the features. Increasing the β value means that minimizing the information flow of the robust features while maximizing the information flow of the non-robust features. Therefore, we will compare classification accuracy for the robust and non-robust features according to the β value. Fig. 9 shows the results when predicting with distilled robust or non-robust features. In the figure, the blue line represents the accuracy of predictions using only the robust features. Also, the red line represents the accuracy of predictions using only the non-robust features. As shown in the figure, as the β value increases, the accuracy of predicting with the robust feature decreases. Since the amount of information flowing into the robust feature decreases as the β value increases, the prediction accuracy with robust features decreases. On the other hand, as the β value increases, the accuracy of predicting with non-robust features increases since the amount of information flowing into the non-robust feature increase. Through the experiment, we demonstrated how the information bottleneck affects the information flow of the robust and non-robust features.

VI. DISCUSSION

In this research, we generate robust proxies that have robust feature representations. Then, we train the model so that the features resemble the representation of those proxies. For future direction, there are many ways that could improve the proposed method. For example, it is possible to train a network that only leverages robust features by manually masking the non-robust features during the training. Also, exploiting multiple proxies from the same class or generating a more reliable proxy could be a good future direction.

VII. CONCLUSION

In this paper, we introduce the intriguing, yet not explored aspect of adversarial training that explicitly learns adversarially robust features. Many works have demonstrated that adversarial vulnerability mainly stems from the non-robust components of learned features, while how to explicitly learn robust features is not explored. To tackle the problem, we manually generate adversarially robust features and propose a novel training framework called robust proxy learning that explicitly learns robust features. To this end, through the CEO algorithm, we generate class representative robust features called robust proxies. During the training, DNNs explicitly learn the representation of robust proxies through the proposed robust proxy learning framework. For each proxy, we pull the data of the same class close to the proxy and push others away in the feature space, allowing the model to explicitly learn adversarially robust features. Extensive experimental results suggest that the proposed method can improve the robustness of existing AT methods under stronger attacks and be general and flexible enough to be adopted on any AT methods. We believe that the proposed method could shed new insight into utilizing robust perturbation for adversarial robustness.

REFERENCES

- A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [2] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates *et al.*, "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.
- [3] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *Advances in neural information processing systems*, vol. 26, 2013.
- [4] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," arXiv preprint arXiv:1412.6572, 2014.
- [5] A. Nguyen, J. Yosinski, and J. Clune, "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 427–436.
- [6] 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.
- [7] 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.
- [8] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *Proceedings of the 10th ACM* workshop on artificial intelligence and security, 2017, pp. 15–26.
- [9] 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.
- [10] F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1778–1787.
- [11] G. S. Dhillon, K. Azizzadenesheli, Z. C. Lipton, J. Bernstein, J. Kossaifi, A. Khanna, and A. Anandkumar, "Stochastic activation pruning for robust adversarial defense," arXiv preprint arXiv:1803.01442, 2018.
- [12] C. Guo, M. Rana, M. Cisse, and L. Van Der Maaten, "Countering adversarial images using input transformations," arXiv preprint arXiv:1711.00117, 2017.

- [13] T. Pang, K. Xu, C. Du, N. Chen, and J. Zhu, "Improving adversarial robustness via promoting ensemble diversity," in *International Conference* on Machine Learning. PMLR, 2019, pp. 4970–4979.
- [14] B.-K. Lee, J. Kim, and Y. M. Ro, "Masking adversarial damage: Finding adversarial saliency for robust and sparse network," *arXiv preprint* arXiv:2204.02738, 2022.
- [15] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu, "Improving adversarial robustness requires revisiting misclassified examples," in *International Conference on Learning Representations*, 2019.
- [16] 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.
- [17] R. Rade and S.-M. Moosavi-Dezfooli, "Reducing excessive margin to achieve a better accuracy vs. robustness trade-off," in *International Conference on Learning Representations*, 2022. [Online]. Available: https://openreview.net/forum?id=Azh9QBQ4tR7
- [18] T. Bai, J. Luo, J. Zhao, B. Wen, and Q. Wang, "Recent advances in adversarial training for adversarial robustness," *arXiv preprint* arXiv:2102.01356, 2021.
- [19] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," in *International conference on machine learning*. PMLR, 2018, pp. 274–283.
- [20] C. Mao, Z. Zhong, J. Yang, C. Vondrick, and B. Ray, "Metric learning for adversarial robustness," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [21] L. Fan, S. Liu, P.-Y. Chen, G. Zhang, and C. Gan, "When does contrastive learning preserve adversarial robustness from pretraining to finetuning?" *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [22] Y. Zhong and W. Deng, "Adversarial learning with margin-based triplet embedding regularization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6549–6558.
- [23] M. Kim, J. Tack, and S. J. Hwang, "Adversarial self-supervised contrastive learning," Advances in Neural Information Processing Systems, vol. 33, pp. 2983–2994, 2020.
- [24] Z. Jiang, T. Chen, T. Chen, and Z. Wang, "Robust pre-training by adversarial contrastive learning," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 16199– 16210. [Online]. Available: https://proceedings.neurips.cc/paper/2020/ file/ba7e36c43aff315c00ec2b8625e3b719-Paper.pdf
- [25] S. Gowal, P.-S. Huang, A. van den Oord, T. Mann, and P. Kohli, "Selfsupervised adversarial robustness for the low-label, high-data regime," in *International Conference on Learning Representations*, 2021. [Online]. Available: https://openreview.net/forum?id=bgQek2O63w
- [26] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," *Advances in neural information processing systems*, vol. 32, 2019.
- [27] J. Kim, B.-K. Lee, and Y. M. Ro, "Distilling robust and non-robust features in adversarial examples by information bottleneck," *Advances* in Neural Information Processing Systems, vol. 34, 2021.
- [28] K. Roth, Y. Kilcher, and T. Hofmann, "The odds are odd: A statistical test for detecting adversarial examples," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5498–5507.
- [29] A. Shafahi, W. R. Huang, C. Studer, S. Feizi, and T. Goldstein, "Are adversarial examples inevitable?" in *International Conference* on Learning Representations, 2019. [Online]. Available: https: //openreview.net/forum?id=r11WUoA9FQ
- [30] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness may be at odds with accuracy," *arXiv preprint* arXiv:1805.12152, 2018.
- [31] S. Yang, T. Guo, Y. Wang, and C. Xu, "Adversarial robustness through disentangled representations," in *Proceedings of the AAAI Conference* on Artificial Intelligence, vol. 35, no. 4, 2021, pp. 3145–3153.
- [32] J. Kim, B.-K. Lee, and Y. M. Ro, "Demystifying causal features on adversarial examples and causal inoculation for robust network by adversarial instrumental variable regression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), June 2023, pp. 12 302–12 312.
- [33] J. Zhang, J. Zhu, G. Niu, B. Han, M. Sugiyama, and M. Kankanhalli, "Geometry-aware instance-reweighted adversarial training," arXiv preprint arXiv:2010.01736, 2020.
- [34] D. Jakubovitz and R. Giryes, "Improving dnn robustness to adversarial attacks using jacobian regularization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 514–529.

- [35] Y. Zhong and W. Deng, "Adversarial learning with margin-based triplet embedding regularization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [36] H. Wang, Y. Deng, S. Yoo, H. Ling, and Y. Lin, "Agkd-bml: Defense against adversarial attack by attention guided knowledge distillation and bi-directional metric learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 7658–7667.
- [37] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," in *International Conference on Learning Representations*, 2017. [Online]. Available: https://openreview. net/forum?id=HyxQzBceg
- [38] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [39] Y. Le and X. Yang, "Tiny imagenet visual recognition challenge," CS 231N, vol. 7, no. 7, p. 3, 2015.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision* and pattern recognition, 2016, pp. 770–778.
- [41] S. Zagoruyko and N. Komodakis, "Wide residual networks," arXiv preprint arXiv:1605.07146, 2016.
- [42] F. Croce and M. Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," in *International conference on machine learning*. PMLR, 2020, pp. 2206–2216.
- [43] —, "Minimally distorted adversarial examples with a fast adaptive boundary attack," in *International Conference on Machine Learning*. PMLR, 2020, pp. 2196–2205.
- [44] M. Andriushchenko, F. Croce, N. Flammarion, and M. Hein, "Square attack: a query-efficient black-box adversarial attack via random search," in *European Conference on Computer Vision*. Springer, 2020, pp. 484– 501.
- [45] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. Goodfellow, A. Madry, and A. Kurakin, "On evaluating adversarial robustness," arXiv preprint arXiv:1902.06705, 2019.
- [46] F. Tramer, N. Carlini, W. Brendel, and A. Madry, "On adaptive attacks to adversarial example defenses," *Advances in neural information processing systems*, vol. 33, pp. 1633–1645, 2020.
- [47] F. Tramer, "Detecting adversarial examples is (nearly) as hard as classifying them," in *International Conference on Machine Learning*. PMLR, 2022, pp. 21 692–21 702.
- [48] B. C. Kim, J. U. Kim, H. Lee, and Y. M. Ro, "Revisiting role of autoencoders in adversarial settings," in 2020 IEEE International Conference on Image Processing (ICIP). IEEE, 2020, pp. 1856–1860.
- [49] M. Zhou and V. M. Patel, "Enhancing adversarial robustness for deep metric learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15325–15334.



YONG MAN RO (Senior Member, IEEE) is a Professor of Electrical Engineering and the Director of the Center for Applied Research in Artificial Intelligence (CARAI) at the Korea Advanced Institute of Science and Technology (KAIST). He received his B.S. degree from Yonsei University, Seoul, South Korea, and his M.S. and Ph.D. degrees from KAIST. He has been a Researcher at Columbia University, a Visiting Researcher at the University of California at Irvine, and a Research Fellow of the University of California at Berkeley. He was

also a Visiting Professor with the Department of Electrical and Computer Engineering, University of Toronto, Canada. His research interests span image and video systems, including image processing, computer vision, multimodal learning, vision-language learning, and object detection. He received the Young Investigator Finalist Award of ISMRM in 1992 and the Year's Scientist Award (Korea) in 2003. He is an Associate Editor for IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY and has served as an Associate Editor for IEEE SIGNAL PROCESSING LETTERS. He has also served as a TPC member, program chair, and special session organizer for many international conferences.



HONG JOO LEE received the B.S. degree from Ajou University, Suwon, South Korea, in 2016, and the M.S. and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2018 and 2023. His research interests include deep learning, machine learning, medical image segmentation, and adversarial robustness.