

Defensive Few-shot Adversarial Learning

Wenbin Li¹, Lei Wang², Xingxing Zhang³, Jing Huo¹, Yang Gao¹, Jiebo Luo⁴
¹Nanjing University, China, ²University of Wollongong, Australia
³Beijing Jiaotong University, China, ⁴University of Rochester, USA

Abstract

The robustness of deep learning models against adversarial attacks has received increasing attention in recent years. However, both deep learning and adversarial training rely on the availability of a large amount of labeled data and usually do not generalize well to new, unseen classes when only a few training samples are accessible. To address this problem, we explicitly introduce a new challenging problem – how to learn a robust deep model with limited training samples per class, called defensive few-shot learning in this paper. Simply employing the existing adversarial training techniques in the literature cannot solve this problem. This is because few-shot learning needs to learn transferable knowledge from disjoint auxiliary data, and thus it is invalid to assume the sample-level distribution consistency between the training and test sets as commonly assumed in existing adversarial training techniques. In this paper, instead of assuming such a distribution consistency, we propose to make this assumption at a task-level in the episodic training paradigm in order to better transfer the defense knowledge. Furthermore, inside each task, we design a task-conditioned distribution constraint to narrow the distribution gap between clean and adversarial examples at a sample-level. These give rise to a novel mechanism called multi-level distribution based adversarial training (MDAT) for learning transferable adversarial defense. In addition, a unified \mathcal{F}_β score is introduced to evaluate different defense methods under the same principle. Extensive experiments demonstrate that MDAT achieves higher effectiveness and robustness over existing alternatives in the few-shot case.

1. Introduction

Deep convolutional neural networks (CNNs) [11, 9] have obtained impressive successes on a variety of computer vision tasks especially in image classification [5, 4]. Unfortunately, several pieces of recent work [3, 2] have shown that these CNN models are vulnerable to adversarial examples, which are crafted based on original clean inputs with imperceptible perturbations. Therefore, how to learn

robust models that can effectively defend against various kinds of attacks becomes a crucial problem. Like an offense and defense game, recent adversarial defense studies continue to develop new attack techniques [1, 23]. Accordingly, more effective defense mechanisms that are robust to these attacks have also been developed [13, 14]. This has considerably improved the robustness of deep learning models [26, 6, 21].

However, both traditional and adversarially trained deep models rely on a large amount of labeled data available for each class. In many real applications, we may only have access to a few labeled samples for new, unseen classes (*e.g.*, rare species or medical diseases). Learning a classifier when only a few training samples are available for each class is well known as few-shot learning in the literature. In the few-shot setting, deep models become more vulnerable to adversarial attacks due to the serious lack of training samples [18]. How to learn an effective and robust model with limited training data becomes a new challenging problem, and we call it *defensive few-shot learning* in this paper.

1.1. Four Questions

In this paper, we concentrate on defensive few-shot image classifications and attempt to answer four key questions around this new problem: (1) what is the main difference between generic image classification and defensive few-shot image classification? (2) how to transfer adversarial defense knowledge from one sample distribution to another? (3) how to narrow the distribution gap between clean and adversarial examples in the few-shot setting? (4) last but not least, how to compare different defense methods in a convenient and fair manner?

For the first question, a recent important work [18] proved that adversarially robust generalization requires the access to more data. However, the few-shot setting has access to significantly fewer training samples than generic image classification. This makes defensive few-shot image classification more difficult to achieve than the robust generic image classification. In addition, generic image classification can usually safely assume the sample-level distribution consistency between the training and test sets.

However, this assumption cannot be made any more in the few-shot case. This is because due to the serious lack of training samples, few-shot learning usually has to resort to a large-scale but class-disjoint auxiliary set to learn transferable knowledge. It means that the actual training (*i.e.*, auxiliary) set often has a somewhat different sample distribution from the test set of the target few-shot task. Such a difference naturally leads to the second question, *i.e.*, how to transfer adversarial defense between two different sample distributions, which has not been sufficiently considered in the literature. The distribution discrepancy is also the reason why we cannot directly apply the existing adversarial defense techniques to address the defensive few-shot learning problem.

As for the third question, it is well known that there is usually a large sample distribution gap between clean and adversarial examples [19]. This distribution gap plays a key role in making deep models vulnerable to adversarial examples. Note that such a gap does not exist in traditional few-shot learning which has not considered adversarial attacks. However, in the new problem of defensive few-shot learning introduced in this paper, we have to tackle this distribution gap in order to improve the robustness of classification models. For the last question, we find that the existing adversarial defense methods often report two kinds of classification accuracy, *i.e.*, clean example accuracy and adversarial example accuracy, to show the effectiveness of their proposed defense methods. This is understandable because a method better at defending against adversarial examples could be worse at classifying clean examples, or vice versa. Nevertheless, this makes a direct comparison of different methods awkward. Therefore, it is desirable to have a unified criterion to facilitate the evaluation and comparison of different defensive methods under the same principle, which has been largely overlooked in the existing literature.

1.2. Our Contributions

Two kinds of distribution gaps (*i.e.*, our challenges) exist and entangle in *defensive few-shot learning*, making it different from both generic adversarial training and few-shot learning. The first is the distribution gap between the training (auxiliary) set and test set, driving us to think about how to transfer adversarial defense in this case. The second is the distribution gap between clean and adversarial examples, making us to consider how to narrow this gap along with the first gap.

To address these issues and pursue the answers to the above four questions, we propose a *Multi-level Distribution based Adversarial Training (MDAT)* mechanism for defensive few-shot learning. The main contributions are:

- We analyze and explore how to learn a robust deep model with very limited training examples per class for the first time in the literature.

- We propose an *episode-based adversarial training mechanism* to transfer adversarial defense across different sample distributions by assuming the distribution consistency at the **task-level**.
- We propose a *task-conditioned distribution constraint*, which is adaptative to different few-shot tasks, to narrow the distribution gap between the clean and adversarial examples inside each task at the **sample-level**.
- We introduce a unified \mathcal{F}_β score to conveniently evaluate the overall performance of different defense methods under the same principle.
- We build upon a state-of-the-art few-shot learning method and conduct extensive experiments on two benchmark datasets. This provides rich baseline results for this new problem, *i.e.*, defensive few-shot learning, to facilitate future research on this topic.

2. Preliminaries

In this section, we first briefly present the definition of the new topic, *i.e.*, *defensive few-shot learning*. Next, we introduce some related notations and terminologies used in this new topic, and review two classic attack methods.

2.1. Defensive Few-shot Learning

A few-shot task normally consists of a support set \mathcal{S} and a query set \mathcal{Q} , where \mathcal{S} contains C different classes with K images per class. Given \mathcal{S} , the target of few-shot learning is to infer the correct class label from the C classes for each unlabeled sample in \mathcal{Q} . This setting can be seen as a C -way K -shot classification setting. Generally, the number of K is small (*e.g.*, 1 or 5), making it almost impossible to learn an effective classifier only from the support set \mathcal{S} . To tackle this problem, an additional auxiliary set \mathcal{A} is usually adopted to learn transferable knowledge to help the classification on \mathcal{Q} . Importantly, \mathcal{A} contains a larger number of classes and samples than \mathcal{S} , while it has a totally disjoint class space with \mathcal{S} .

Unlike the generic few-shot learning, here, we mainly focus on how to learn a defensive few-shot classification model to defend against adversarial attacks, *i.e.*, *defensive few-shot learning*. In defensive few-shot learning, we always assume that the adversary is capable of manipulating the query images in the query set \mathcal{Q} , but doesn't have access to the support set \mathcal{S} . Thus, our goal is to learn a robust model which can correctly classify query images no matter if they are manipulated.

2.2. Notation

Let $\Psi(\mathbf{x})$ denote a convolutional neural network based embedding module, which can learn feature representations for any input image \mathbf{x} . We use $f(\Psi(\mathbf{x}), \mathcal{S}) : \mathbb{R}^d \rightarrow \mathbb{R}^C$

as a classifier module to assign a class label y for a query image x within C different classes, according to the support set \mathcal{S} . Denote the true class label of x as y_{true} . Note that these two modules can be integrated into a unified network and trained from scratch in an end-to-end manner. The cost function of $f(\Psi(x), \mathcal{S})$ is denoted as $\mathcal{L}(x, \mathcal{S}, y_{true})$ for simplicity. It is easy to use a small perturbation δ to construct an adversarial image $x^{adv} = x + \delta$ to fool the classifier, making $f(\Psi(x^{adv}), \mathcal{S}) \neq y_{true}$. Generally, the clean image x and the adversarial image x^{adv} are perceptually indistinguishable, which can be bounded by a distance metric $D(x, x^{adv}) \leq \epsilon$. That is to say, ϵ indicates the maximum magnitude of the perturbation. Here, we use ℓ_∞ norm to measure the bound of the perturbation, which means $\|\delta\|_\infty \leq \epsilon$. Note that, all clean images are normalized into a range of $[0, 1]$, and all adversarial images are clipped into the same range. Thus, for a specific perturbation such as $\epsilon = 0.01$, it corresponds to 4 pixels changing in the range of $[0, 255]$. In particular, we employ a white-box attack setting [2] to generate all the training adversarial images as in [10].

2.3. Attack Methods

Two typical attack methods, *i.e.*, FGSM [3] and PGD [13], are usually adopted to generate adversarial examples during the training phase.

Fast Gradient Sign Method (FGSM) [3] is a one-step method, which generates an adversarial example x^{adv} through a single backward propagation of a neural network with respect to the input x .

$$x^{adv} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, \mathcal{S}, y_{true})), \quad (1)$$

where $\nabla_x \mathcal{L}(x, \mathcal{S}, y_{true})$ indicates the gradient of the loss function $\mathcal{L}(x, \mathcal{S}, y_{true})$ with respect to x . FGSM is a commonly used adversary because it is simple, computationally efficient, yet amenable.

Projected Gradient Descent (PGD) [13], a stronger iterative variant of FGSM, is equivalent to *iterative fast gradient sign method (iFGSM)* [10] but with a random starting point. It mainly applies the FGSM iteratively for k times with a small step size α .

$$\begin{aligned} x_0^{adv} &= x + 0.001 \cdot \mathcal{N}(0, I) \\ x_k^{adv} &= \text{Clip}_{x, \epsilon} \{ x_{k-1}^{adv} + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(x, \mathcal{S}, y_{true})) \}, \end{aligned} \quad (2)$$

where $\text{Clip}_{x, \epsilon}(B)$ performs element-wise clipping of B , making each $B_{i,j}$ clipped into the range $[x - \epsilon, x + \epsilon]$. To choose a random starting point, we add a random Gaussian noise $\mathcal{N}(0, I)$ into the input x like [26, 25], where $\mathcal{N}(0, I)$ denotes a Gaussian distribution with zero mean and identity variance. Following [10], we set $\alpha = \frac{1}{255}$, and the number of iterations $k = \min(255\epsilon + 4, 255\epsilon \cdot 1.25)$ in the $[0, 1]$ space. PGD has a stronger attack strength than FGSM,

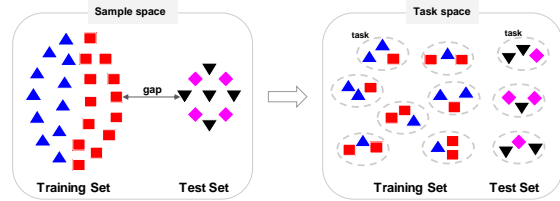


Figure 1. Converting the distribution consistency assumption in the sample level into the task level. Each geometric shape indicates one sample and each color denotes one class.

and thus the networks trained with PGD adversaries have stronger robustness against a wide range of other first-order attacks [13].

3. MDAT: Multi-level Distribution based Adversarial Training

3.1. Task-level Distribution

In generic image classification, we usually take a distribution consistency assumption between the training and test sets (*i.e.*, independently and identically distributed data), guaranteeing the model trained on the training set can generalize to the test set. However, because of the lack of training samples in the support set \mathcal{S} , few-shot learning usually resorts to a class-disjoint auxiliary (training) set \mathcal{A} to learn transferable knowledge for \mathcal{S} . Since the sample distribution of \mathcal{A} is relatively different from that of \mathcal{S} , the generalization performance on the target data sets (*i.e.*, \mathcal{S} and \mathcal{Q}) cannot be well guaranteed, so does the generalization of adversarial training. As an illustration, from the left side of Fig. 1, we can observe that there may be a large distribution gap between training and test sets in the sample space because these two sets are class-disjoint, making the adversarial defense hardly transferable.

Fortunately, we can assume the above distribution consistency in a task space instead of the sample space (see the right side of Fig. 1). The task space is composed of a lot of similar tasks, *e.g.*, few-shot tasks. From the perspective of task-level distribution, we can construct a large number of adversarial-based few-shot tasks within the auxiliary set \mathcal{A} , by simulating the target adversarial-based few-shot task. In doing so, the sample distribution gap can be dealt with by taking the task distribution consistency assumption.

In fact, this kind of task construction strategy has been adopted in the *episodic training mechanism* [22] in traditional few-shot learning, demonstrating very promising performance. However, the effectiveness of this mechanism under the adversarial-based few-shot tasks has not been investigated. In this work, we reinterpret this mechanism from the perspective of task distribution consistency and develop transferable adversarial defense upon it. To the best of our knowledge, addressing this task in such a setting is

the first time in the literature.

Episode-based Adversarial Training. According to the above task-level distribution consistency assumption, we propose the following *episode-based adversarial training*. Specifically, let $\{\langle \mathcal{S}_1, \mathcal{Q}_1 \rangle, \dots, \langle \mathcal{S}_n, \mathcal{Q}_n \rangle\}$ be a set of few-shot tasks randomly sampled from the auxiliary set \mathcal{A} , the objective function of our episode-based adversarial training can be formulated as follows:

$$\Gamma = \arg \min_{\theta} \sum_{i=1}^n \sum_{\mathbf{x} \in \mathcal{Q}_i} \left(\mathcal{L}(\mathbf{x}, \mathcal{S}_i, y_{true}) + \max_{\mathbf{x}^{adv} = \mathbf{x} + \delta} \mathcal{L}(\mathbf{x}^{adv}, \mathcal{S}_i, y_{true}) \right), \quad (3)$$

where θ denotes the parameter set of the defensive few-shot classification model. The core idea here is to simulate the target adversarial-based few-shot task by conducting a lot of similar few-shot tasks with the auxiliary training set. In this way, we can build a task-based space to mitigate the gap among different sample distributions. At each training step, we generate adversarial examples inside each of the sampled tasks according to the current model, and meanwhile, inject these adversarial examples into the training set. Both clean and adversarial examples will be used to train this model, enhancing its capability to defend adversarial attacks.

3.2. Sample-level Distribution

Besides the above mentioned sample distribution gap between the training (auxiliary) and test sets, there is also a significant sample distribution gap between the clean and adversarial examples inside each few-shot task. The traditional few-shot learning cannot deal with such a gap because it has not considered adversarial attacks. On the other hand, the generic adversarial training only assigns the same ground truth to both clean and the corresponding adversarial samples (represented as global feature vectors), while ignoring the underlying distributions between these two.

To tackle this issue, we propose a task-conditioned and distribution-based regularizer \mathcal{L}_{reg} inside each few-shot task $\langle \mathcal{S}_i, \mathcal{Q}_i \rangle$ used for training. This regularizer effectively integrates a *local-descriptor distribution consistency* constraint \mathcal{L}_{fea} with a *class-prediction consistency* constraint \mathcal{L}_{class} , which are formulated as follows and explained in further

$$\mathcal{L}_{reg}(\mathbf{x}, \mathbf{x}^{adv}, \mathcal{S}_i) = \mathcal{L}_{fea}(\Psi(\mathbf{x}), \Psi(\mathbf{x}^{adv})) + \mathcal{L}_{class}(f(\Psi(\mathbf{x}), \mathcal{S}_i), f(\Psi(\mathbf{x}^{adv}), \mathcal{S}_i)), \quad (4)$$

where $\mathbf{x} \in \mathcal{Q}_i$, $\Psi(\cdot)$ denotes the local descriptor based embedding module, and $f(\cdot, \cdot)$ indicates the classifier module.

Local-descriptor distribution consistency. \mathcal{L}_{fea} is designed to narrow the sample distribution gap between the clean and adversarial examples in a given task. Note that this is conducted in the local feature space as inspired by

the recent state-of-the-art few-shot learning method [12]. It represents each image as a set of deep local descriptors obtained in a 3D convolutional feature map instead of a pooled global feature vector. Suppose the local descriptors of clean and adversarial examples follow multivariate normal distributions, *i.e.*, $\Psi(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\Psi(\mathbf{x}^{adv}) \sim \mathcal{N}(\boldsymbol{\mu}^{adv}, \boldsymbol{\Sigma}^{adv})$, we design a novel task-conditioned distribution (TCD) measure as follows:

$$\mathcal{L}_{fea} = (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv})^\top \boldsymbol{\Sigma}_{\mathcal{S}_i}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv}) + \text{Tr} \left(\left\| \boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\Sigma} \boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}} - \boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\Sigma}^{adv} \boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}} \right\|_F^2 \right), \quad (5)$$

where $\text{Tr}(\cdot)$ is the trace operation and $\boldsymbol{\Sigma}_{\mathcal{S}_i}^{-1}$ denotes the inverse covariance matrix of the support set \mathcal{S}_i in the current few-shot task. Specifically, we use the local descriptors of all the samples of all the classes in \mathcal{S}_i to calculate one single covariance matrix $\boldsymbol{\Sigma}_{\mathcal{S}_i}$, which can be seen as a meta-feature of the current few-shot task. Thus, Eq.(5) is conditioned on the current few-shot task by transforming both clean and adversarial examples into a common space based on the support set \mathcal{S}_i . In particular, the first term of Eq.(5) is a squared Mahalanobis distance between $\boldsymbol{\mu}$ and $\boldsymbol{\mu}^{adv}$, depending on \mathcal{S}_i . The second term aims to measure the distribution distance between the clean and adversarial examples with the second-order information, which is also depended on \mathcal{S}_i . The purpose of $\boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}}$ is to transform both clean and adversarial examples of the query set \mathcal{Q}_i into the same space defined by the current task (*see the appendix for more details*). However, the complicated matrix square root (*i.e.*, $\boldsymbol{\Sigma}_{\mathcal{S}_i}^{-\frac{1}{2}}$) makes the calculation of Eq.(5) time-consuming. Fortunately, according to Theorem 1, we can convert Eq.(5) into an efficient form without matrix square root to accelerate the calculation.

It is also worth mentioning that it is difficult to use pooled global features to effectively represent the sample distribution (say, computing the covariance matrix) of both clean and adversarial examples, due to the limited examples under few-shot setting. Instead, as aforementioned, we employ the rich local descriptors of the convolutional layer to represent each image and then measure the distribution consistency between clean and adversarial examples in the local feature space.

Theorem 1. *Suppose $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2$ and $\boldsymbol{\Sigma}$ are all positive semi-definite matrices, and $\boldsymbol{\Sigma}^{-\frac{1}{2}}$ is the square root of the inverse of $\boldsymbol{\Sigma}$, we have*

$$\begin{aligned} \text{Tr} \left(\left\| \boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}^{-\frac{1}{2}} - \boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}^{-\frac{1}{2}} \right\|_F^2 \right) = \\ \text{Tr} [\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}^{-1} \cdot \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}^{-1}] - \\ 2 \text{Tr} [\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}^{-1} \cdot \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}^{-1}] + \\ \text{Tr} [\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}^{-1} \cdot \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}^{-1}]. \end{aligned} \quad (6)$$

Proof. Let $\hat{\boldsymbol{\Sigma}}_1 = \boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}^{-\frac{1}{2}}$ and $\hat{\boldsymbol{\Sigma}}_2 = \boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}^{-\frac{1}{2}}$,

we have

$$\begin{aligned}
& \text{Tr} \left(\|\Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-\frac{1}{2}} - \Sigma^{-\frac{1}{2}} \Sigma_2 \Sigma^{-\frac{1}{2}}\|_F^2 \right) \\
&= \text{Tr} \left(\|\hat{\Sigma}_1 - \hat{\Sigma}_2\|_F^2 \right) \\
&= \text{Tr} \left[(\hat{\Sigma}_1 - \hat{\Sigma}_2)^\top (\hat{\Sigma}_1 - \hat{\Sigma}_2) \right] \\
&= \text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_1 - 2 \hat{\Sigma}_1^\top \hat{\Sigma}_2 + \hat{\Sigma}_2^\top \hat{\Sigma}_2 \right] \\
&= \text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_1 \right] - 2 \text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_2 \right] + \text{Tr} \left[\hat{\Sigma}_2^\top \hat{\Sigma}_2 \right].
\end{aligned} \tag{7}$$

Specifically, according to the basic properties of matrix, for $\text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_1 \right]$, we can obtain

$$\begin{aligned}
\text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_1 \right] &= \text{Tr} \left[(\Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-\frac{1}{2}})^\top (\Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-\frac{1}{2}}) \right] \\
&\quad (\text{both } \Sigma_1 \text{ and } \Sigma^{-\frac{1}{2}} \text{ are symmetric}) \\
&= \text{Tr} \left[\Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-\frac{1}{2}} \right] \\
&\quad (\text{because } \Sigma^{-1} = \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}}) \\
&= \text{Tr} \left[\Sigma^{-\frac{1}{2}} \Sigma_1 \Sigma^{-1} \Sigma_1 \Sigma^{-\frac{1}{2}} \right] \\
&\quad (\text{because } \text{Tr}(ABC) = \text{Tr}(CBA)) \\
&= \text{Tr} \left[\Sigma_1 \Sigma^{-1} \cdot \Sigma_1 \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} \right] \\
&= \text{Tr} \left[\Sigma_1 \Sigma^{-1} \cdot \Sigma_1 \Sigma^{-1} \right].
\end{aligned} \tag{8}$$

Similarly, we can easily obtain

$$\begin{aligned}
\text{Tr} \left[\hat{\Sigma}_1^\top \hat{\Sigma}_2 \right] &= \text{Tr} \left[\Sigma_1 \Sigma^{-1} \cdot \Sigma_2 \Sigma^{-1} \right] \\
\text{Tr} \left[\hat{\Sigma}_2^\top \hat{\Sigma}_2 \right] &= \text{Tr} \left[\Sigma_2 \Sigma^{-1} \cdot \Sigma_2 \Sigma^{-1} \right].
\end{aligned} \tag{9}$$

Bring Eq.(8) and Eq.(9) back to Eq.(7), we can obtain Eq.(6). Hence proved. \square

Class-prediction consistency. As for \mathcal{L}_{class} , the goal is to make the class predictions of both clean and adversarial examples similar. In other words, the probability distributions of the predicted classes of clean and adversarial examples should be similar. For this end, we employ a simple ℓ_2 norm to conduct this loss, that is,

$$\mathcal{L}_{class} = \|f(\Psi(\mathbf{x}), \mathcal{S}_i) - f(\Psi(\mathbf{x}^{adv}), \mathcal{S}_i)\|_2^2. \tag{10}$$

Thus, the overall optimization formulation of the proposed MDAT framework is

$$\begin{aligned}
\Gamma = \arg \min_{\theta} & \sum_{i=1}^n \sum_{\mathbf{x} \in \mathcal{Q}_i} \left(\mathcal{L}(\mathbf{x}, \mathcal{S}_i, y_{true}) + \right. \\
& \left. \max_{\mathbf{x}^{adv} = \mathbf{x} + \delta} \mathcal{L}(\mathbf{x}^{adv}, \mathcal{S}_i, y_{true}) + \lambda \cdot \mathcal{L}_{reg}(\mathbf{x}, \mathbf{x}^{adv}, \mathcal{S}_i) \right),
\end{aligned} \tag{11}$$

where λ is a balancing parameter and $\mathcal{L}(\cdot)$ denotes the cross-entropy loss.

3.3. A Unified Evaluation Criterion

Training a model with adversarial examples can indeed improve the robustness of this model. At the same time, it could jeopardise the performance of this model on the original clean examples. There may exist a trade-off between

robustness (against adversarial examples) and accuracy (on clean examples). Several recent works have been trying to show this issue [21, 25]. In the literature, adversarial training based work generally reports two kinds of classification accuracy, *i.e.*, accuracy on the clean examples and accuracy on the adversarial examples. However, in real cases, it could be awkward to compare two defense methods overall with two accuracies. To this end, inspired by the case of Recall and Precision in information retrieval, we introduce \mathcal{F}_β score as a unified criterion to evaluate different defense methods under the same principle. To be specific, \mathcal{F}_β is formulated with both clean accuracy ACC_{clean} and adversarial accuracy ACC_{adv} as

$$\mathcal{F}_\beta = (1 + \beta^2) \cdot \frac{ACC_{clean} \cdot ACC_{adv}}{\beta^2 \cdot ACC_{clean} + ACC_{adv}}. \tag{12}$$

For instance, if we would like to maintain a high accuracy on the original clean examples, and meanwhile, improve the robustness as higher as possible, we can use $\mathcal{F}_{0.5}$. In contrast, if we mainly concern the robustness, \mathcal{F}_2 will be an alternative option. In addition, we can also generate the curve of \mathcal{F}_β by varying β as a more comprehensive way to compare different methods (see the appendix).

4. Experiments

In this section, we perform few-shot image classification on two benchmark datasets by applying our proposed MDAT mechanism to learn defensive few-shot classifiers. The main targets include: (1) is the task-level distribution consistency assumption effective? (2) can the task-conditioned regularization improve the robustness and accuracy?

Datasets. Two datasets are picked as the benchmarks, *i.e.*, *miniImageNet* [22] and CIFAR-100 [8]. *miniImageNet* consists of 100 classes, and there are 600 images in each class with a resolution of 84×84 . We follow [17] and take 64, 16 and 20 classes for training (auxiliary), validation and test, respectively. As for CIFAR-100, we follow [20] and take 60, 20 and 20 classes for training (auxiliary), validation and test.

Network Architecture. A commonly used four-layer convolutional neural network [24, 12] in few-shot learning is adopted as the embedding module. It consists of four convolutional blocks, each of which contains a convolutional layer, a batch normalization layer and a LeakyReLU layer. Note that other networks (*e.g.*, ResNet [5]) can be used as alternative, which is not the main concern of this paper.

Attack Setting. We apply FGSM [3] and PGD [13] attackers to find adversarial examples for training two kinds of robust models, *i.e.*, FGSM-based and PGD-based, respectively. During training, we follow [10] and randomly choose a perturbation ϵ for each training few-shot task from a normal distribution in the range of $[0, 0.01]$. Moreover,

Table 1. Comparison (%) of DN4 and DN4-AT with or without episodic training on *miniImageNet*. A FGSM attacker is adopted.

Method	Episode	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
			ACC	$\mathcal{F}_{0.5}$	ACC	$\mathcal{F}_{0.5}$	ACC	$\mathcal{F}_{0.5}$
DN4	without	53.64	22.03	41.67	18.30	38.69	16.43	36.91
DN4-AT	without	52.59	44.59	50.76	35.15	47.84	31.93	46.56
DN4	with	70.84	17.25	43.69	7.23	25.67	5.57	21.18
DN4-AT	with	67.30	55.23	64.48	41.91	60.02	32.89	55.65

for each input clean image (*i.e.*, query image), we construct its corresponding adversarial counterpart. The hyperparameter λ in Eq.(11) is empirically set as 0.5 for all the experiments. During test, we evaluate the robustness of the trained models in defense of three levels of attacks, *i.e.*, $\epsilon = \{0.003, 0.007, 0.01\}$.

Experimental Setting. For fairness and simplicity, all experiments are conducted around a 5-way 5-shot task on the benchmark datasets based on one state-of-the-art few-shot learning method, *i.e.*, DN4 [12]. The reason why we choose DN4 is that it is one of the latest methods and uses local descriptors. Meanwhile, we highlight that in principle, our proposed MDAT mechanism is general and can be applied into any other few-shot learning method.

In the training stage, we employ the episode-based adversarial training to train all the models, following the assumption of task-level distribution consistency. Specifically, we train all the models for 40 epochs, and in each epoch, we randomly sample and construct 10,000 few-shot tasks (*i.e.*, episodes). In each conducted few-shot task, there are 5 support classes with 5 support images and 10 query images per class. Adam algorithm [7] is used to update all the models. The initial learning rate is set as 0.005 and cut in half per 10 epochs. Importantly, all the models are trained from scratch in an end-to-end manner. In the test stage, 3000 few-shot tasks are constructed from the test set for evaluation. The top-1 mean accuracy and the introduced \mathcal{F}_β score are taken as the evaluation criteria. In particular, we set $\beta = 0.5$ in our experiments, which means that the clean accuracy (the weight is 4/5) is regarded as more important than the adversarial accuracy (the weight is 1/5). This considers the fact that clean samples are more common in practice.

4.1. Adversarial Defense Transfer

Our main concern is how to transfer the adversarial defense from a sample distribution to another different one. We have a hypothesis that the task-level distribution consistency assumption is able to improve the generalization performance of defense ability of our model. To verify this hypothesis, we conduct a comparison experiment on *miniImageNet*.

We first use the standard adversarial training on the auxiliary set \mathcal{A} (64 classes), and then transfer the adversarially trained model on the test set (20 classes). Specifically,

we use the standard adversarial training [3] to train a 64-classes classification network on \mathcal{A} , and directly perform DN4 without additional training (because the classification module of DN4 is non-parametric) on the test set based on this pre-trained network. As for the comparison method, we train DN4 on the auxiliary set \mathcal{A} with the proposed episode-based adversarial training (*i.e.*, Eq.(3)), and evaluate it on the test set. Furthermore, two non-adversarially trained DN4s are additionally taken as baselines.

All the results are reported in Table 1. As seen, without adversarial training, DN4 is much vulnerable to adversarial attacks, dropping its accuracy from 53.64% (clean accuracy) to 16.43% (adversarial accuracy) when $\epsilon = 0.01$. In contrast, the adversarially trained DN4-AT can indeed defend the adversarial attacks as expected (from 16.43% to 31.93%). However, both clean and adversarial accuracies of DN4-AT are still far from the normal accuracy on the *miniImageNet* dataset (70.84%).

Fortunately, the episodic training makes DN4-AT perform much better than the non-episode based DN4-AT on both clean and adversarial examples. For instance, when $\epsilon = 0.003$, the episode-based DN4-AT obtains 10.64% and 13.72% improvements over the non-episode based one on the adversarial accuracy and $\mathcal{F}_{0.5}$ score, respectively. More importantly, the clean accuracy is improved from 52.59% to 67.30%, which is much close to the normal clean accuracy (70.84%). Therefore, it means that the episode-based adversarial training can not only transfer the adversarial defense knowledge but also preserve the clean classification knowledge. This verifies that the traditional sample distribution consistency assumption may not guarantee the model’s generalization on both clean and adversarial examples, while the task distribution consistency assumption during episodic training can properly make it. This achieves the first target in Section 4.

4.2. Adversarial Defense on *miniImageNet*

We have verified that the episode-based adversarial training can facilitate the adversarial defense transfer and improve the generalization ability of few-shot models. However, the performance gap between the clean and adversarial accuracies is still large (see Tabel 1). It is necessary to introduce defense techniques to narrow this performance gap.

To further verify the effectiveness of our proposed technique on this challenge, we consider episode-based ad-

Table 2. Comparison (%) of different defense methods on *miniImageNet*. Both training and test are based on a FGSM attacker.

Method	Episode	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
			<i>ACC</i>	$\mathcal{F}_{0.5}$	<i>ACC</i>	$\mathcal{F}_{0.5}$	<i>ACC</i>	$\mathcal{F}_{0.5}$
NT	with	70.84	17.25	43.69	7.23	25.67	5.57	21.18
AT [3]	with	67.30	55.23	64.48	41.91	60.02	32.89	55.65
ATDA [19]	with	67.15	54.18	64.08	39.76	59.01	31.14	54.53
ATDA- ℓ_2	with	67.46	55.11	64.56	40.50	59.53	32.94	55.77
ALP [6]	with	66.31	56.51	64.08	44.12	60.24	36.96	57.22
MDAT (ours)	with	67.27	56.97	64.92	44.82	61.14	37.34	57.97

Table 3. Comparison (%) of different defense methods on *miniImageNet*. Both training and test are based on a PGD attacker.

Method	Episode	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
			<i>ACC</i>	$\mathcal{F}_{0.5}$	<i>ACC</i>	$\mathcal{F}_{0.5}$	<i>ACC</i>	$\mathcal{F}_{0.5}$
NT	with	70.88	19.33	46.22	2.33	10.29	0.80	3.82
AT [3]	with	66.69	53.99	63.69	38.00	57.94	27.51	51.90
ATDA [19]	with	67.28	54.20	64.18	36.56	57.60	25.42	50.61
ATDA- ℓ_2	with	68.39	54.98	65.00	37.69	58.80	27.09	52.40
ALP [6]	with	66.27	56.00	63.92	41.65	59.26	32.24	54.71
MDAT (ours)	with	67.67	56.49	65.09	41.72	60.18	32.44	55.59

versarial training as the default setting and compare our proposed MDAT with other adversarial defense methods. Specifically, three state-of-the-art methods are employed, including the standard adversarial training (AT) [3], Adversarial Logit Pairing (ALP) [6], and Adversarial Training with Domain Adaptation (ATDA) [19]. We reimplement all these methods based on DN4 for a fair comparison. Also, they are all trained and tested using a FGSM attacker. Note that ATDA [19] proposes two losses, *i.e.*, a unsupervised domain adaptation loss (UDA) and a supervised domain adaptation loss (SDA). SDA cannot be used in the few-shot case (classes are varied) because it needs to maintain a log-its center for each fixed class, we reimplement ATDA only with the UDA loss. Instead of the original global features, local descriptors are adopted in ATDA for more promising performances. Additionally, we compare with one variant of ATDA (named ATDA- ℓ_2), where UDA is based on ℓ_2 norm but not original ℓ_1 norm. To make it more intuitive, the normal trained (NT) DN4 is added as a baseline, which is only trained on the clean examples.

The results are provided in Table 2. Compared with these methods, *our proposed MDAT can not only maintain the accuracy on the clean examples, but also improve the robustness on the adversarial examples*. Notably, the proposed MDAT obtains both the highest adversarial accuracies and $\mathcal{F}_{0.5}$ scores on all three levels of adversarial attacks. It can also be seen that ATDA using ℓ_1 based UDA loss cannot improve and even jeopardize the model robustness over the standard adversarial training (*i.e.*, AT). On the contrary, our implemented ATDA- ℓ_2 performs better than ATDA on the robustness and can even improve the clean accuracy. This means that it is reasonable to employ ℓ_2 norm based distribution constraints (*i.e.*, Eq.(5) and Eq.(10)) to improve the generalization performance. ALP gains a large improvement over AT on the robustness especially against

a stronger attacker. For instance, it obtains 4.07% improvements over AT when $\epsilon = 0.01$. However, ALP loses almost 1% accuracy on the original clean examples.

Besides the FGSM attack, we also adopt a stronger attacker PGD to train and test all these methods (see Table 3). As seen, the proposed MDAT is still superior to all the other methods (see $\mathcal{F}_{0.5}$ score). It is also worth noting that it is truly difficult to obtain great improvements on both clean and adversarial accuracies simultaneously, *i.e.*, a good trade-off claimed in [21, 25]. In this sense, the improvements of the proposed MDAT over other comparison methods have their significance. Thus, the second target in Section 4 is achieved.

4.3. Adversarial Defense on CIFAR-100

To show the consistent effectiveness of the proposed MDAT, we conduct few-shot classification tasks on a more challenging dataset, *i.e.*, CIFAR-100 with lower image resolution and more challenging class splits. All the settings are similar to the ones on *miniImageNet*, and we also employ a FGSM (PGD) attacker to train and test all the methods, respectively. The results are reported in Tables 4 and 5. It can be seen that our proposed MDAT obtains consistently superior results than other competitors on both clean and adversarial accuracies, especially on the $\mathcal{F}_{0.5}$ score. For example, in Table 4, when $\epsilon = 0.01$, MDAT obtains 29.69%, 2.87%, 4.65% and 0.8% improvements over NT, AT, ATDA and ALP in terms of the $\mathcal{F}_{0.5}$ score, respectively.

4.4. Ablation Study

An ablation study is conducted to verify the effectiveness of each part of the proposed MDAT mechanism in Table 6. First, the episode-based adversarial training is adopted as the default setting. We degrade MDAT to the simplest variant AT by removing the proposed \mathcal{L}_{reg} regu-

Table 4. Comparison (%) of different defense methods on CIFAR-100. Both training and test are based on a FGSM attacker.

Method	Episode	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
			\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$
NT	with	54.37	10.14	29.03	4.48	16.84	3.68	14.47
AT [3]	with	54.78	40.90	51.29	27.59	45.76	20.81	41.29
ATDA [19]	with	53.53	38.55	49.66	25.15	43.67	19.30	39.51
ALP [6]	with	54.09	42.74	51.36	30.56	46.87	24.19	43.36
MDAT (ours)	with	55.12	42.45	52.01	30.92	47.65	24.60	44.16

Table 5. Comparison (%) of different defense methods on CIFAR-100. Both training and test are based on a PGD attacker.

Method	Episode	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
			\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$
NT	with	54.64	11.87	31.75	1.90	8.34	0.62	2.96
AT [3]	with	54.76	43.02	51.92	27.20	45.53	20.11	40.72
ATDA [19]	with	53.92	40.88	50.68	24.41	43.42	17.61	38.17
ALP [6]	with	54.71	43.27	51.96	29.26	46.60	22.61	41.10
MDAT (ours)	with	54.87	44.02	52.29	29.95	47.04	22.56	41.13

Table 6. Ablation study of our MDAT on *miniImageNet* based on a FGSM attacker.

Episode	AT	\mathcal{L}_{class}	\mathcal{L}_{fea}	Clean	$\epsilon=0.003$		$\epsilon=0.007$		$\epsilon=0.01$	
					\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$	\mathcal{ACC}	$\mathcal{F}_{0.5}$
✓	✓			67.30	55.23	64.48	41.91	60.02	32.89	55.65
✓	✓	✓		66.31	56.51	64.08	44.12	60.24	36.96	57.22
✓	✓		✓	68.22	55.74	65.29	41.61	60.48	32.91	56.16
✓	✓	✓	✓	67.27	56.97	64.92	44.82	61.14	37.34	57.97

larization. From Table 6, we can observe that \mathcal{L}_{class} can dramatically improve the adversarial accuracies while only lose approximately 1% clean accuracy. Unlike \mathcal{L}_{class} , the designed \mathcal{L}_{fea} can perform excellently on the clean images (*i.e.*, 68.22%) and also slightly improve the adversarial accuracies in most cases (*e.g.*, $\epsilon=0.003$ and $\epsilon=0.01$). Therefore, to maintain both the clean and adversarial accuracies, combining AT with \mathcal{L}_{fea} and \mathcal{L}_{class} will be a better option.

5. Discussions on Related Work

There is a large body of adversarial training based related work [26, 15, 25], focusing on the robustness of generic image classification. However, the robustness and defense transferability in the few-shot case have been less considered by the state-of-the-art approaches. Here, we only discuss the most relevant studies. To improve the robustness of semi-supervised classification, Miyato *et al.* [14] propose a semi-supervised virtual adversarial training method (VAT) by calculating the KL divergence between the predictions on the clean examples and the adversarial examples. Similarly, Kannan *et al.* [6] present a symmetric adversarial logit pairing strategy (ALP), encouraging similar embeddings of the clean and the corresponding adversarial examples. Recently, Song *et al.* [19] introduce domain adaptation into adversarial training (ATDA) to learn domain invariant representations for both clean and adversarial domains. The main differences between our MDAT and these methods are: (1) the above methods simply consider the generic image clas-

sification setting, which is actually related to the sample-level distribution constraint (*i.e.*, \mathcal{L}_{reg}) in our MDAT; (2) the regularization terms used in these methods are fixed and cannot adapt to different tasks, while our regularization term (*i.e.*, Eq.(5)) facilitates addressing the task shift since it is explicitly conditioned on the task; (3) instead of only focusing on the global embeddings and class-prediction distributions as in these three related methods, we additionally consider the local-descriptor distribution consistency in a local feature space. This can be more effective for capturing the distributions of both clean and adversarial examples in the few-shot case.

6. Conclusions

In this paper, we introduce a new challenging problem for the first time, *i.e.*, defensive few-shot learning, aiming to learn robust classification models with limited training samples in each class. The key to this problem is how to transfer the common representations and adversarial defense knowledge. To that end, we propose a *Multi-level Distribution based Adversarial Training (MDAT)* mechanism. Extensive experiments have verified that (1) the task-level distribution consistency assumption can guarantee the generalization performance of our MDAT; (2) narrowing the sample-level distribution gap with the proposed task-conditioned distribution constraint can further improve the robustness and accuracy.

A. Further Interpretation of Local-descriptor Distribution Consistency

The essence of \mathcal{L}_{fea} is to transform both clean and adversarial examples into a common feature space that is built by the support set \mathcal{S}_i in the current few-shot task $\langle \mathcal{S}_i, \mathcal{Q}_i \rangle$. In the common feature space, we use the proposed distribution measure \mathcal{L}_{fea} to narrow the distribution gap between the clean and adversarial examples. In this way, the current few-task information is conditioned on the distribution measure, which can further tackle the task shift issue accompanied by the task distribution consistency assumption.

Degrade the formulation of \mathcal{L}_{fea} by removing $\Sigma_{\mathcal{S}_i}^{-1}$ and $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}$ as below,

$$\mathcal{L}_{fea} = (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv})^\top (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv}) + \text{Tr}(\|\Sigma - \Sigma^{adv}\|_F^2), \quad (13)$$

which can be seen as an approximate 2-Wasserstein distance [16] between the clean example distribution and adversarial example distribution.

Since $\Sigma_{\mathcal{S}_i}$ is the covariance matrix of the support set \mathcal{S}_i , it is a positive definite matrix. Thus, $\Sigma_{\mathcal{S}_i}$ can be mathematically decomposed as $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}$, where $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}$ can be seen as a transformation matrix. For each clean image \mathbf{x} , it can be transformed as $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \mathbf{x}$. Therefore, in the new feature space, the covariance matrix of the transformed clean examples is $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}$. Similarly, we can obtain $\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma^{adv} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}$ for the transformed adversarial examples. Then we can convert Eq.(13) into

$$\begin{aligned} \mathcal{L}_{fea} &= (\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\mu} - \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\mu}^{adv})^\top (\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\mu} - \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \boldsymbol{\mu}^{adv}) + \\ &\quad \text{Tr}(\|\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} - \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma^{adv} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}\|_F^2) \\ &= (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv})^\top \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv}) + \\ &\quad \text{Tr}(\|\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} - \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma^{adv} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}\|_F^2) \\ &= (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv})^\top \Sigma_{\mathcal{S}_i}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}^{adv}) + \\ &\quad \text{Tr}(\|\Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} - \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}} \Sigma^{adv} \Sigma_{\mathcal{S}_i}^{-\frac{1}{2}}\|_F^2). \end{aligned} \quad (14)$$

B. Qualitative comparison: Curve of \mathcal{F}_β Score

In the main paper, we prefer to use $\mathcal{F}_{0.5}$ score to obtain quantitative results for different comparison methods. In fact, we can also generate the curve of \mathcal{F}_β by varying the value of β for qualitative comparison.

For example, we plot the curve of F_β by varying β from 0 to 2 according to the results in Table 2 ($\epsilon = 0.003$) in the main paper. We can see that our MDAT performs consistently superior to other methods for any value of β (see the above Figure 2).

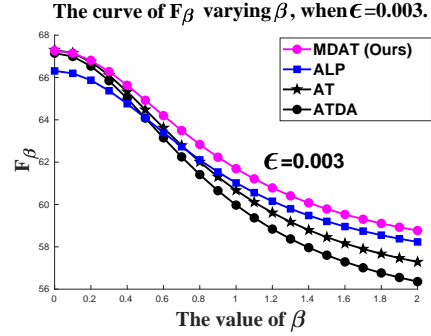


Figure 2. The curves of \mathcal{F}_β of different methods by varying the value of β from 0 to 2.

References

- [1] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, 2018.
- [2] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defenses: A survey. *arXiv*, 2018.
- [3] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv*, 2014.
- [4] Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, Gang Wang, Jianfei Cai, et al. Recent advances in convolutional neural networks. *Pattern Recognition*, 77:354–377, 2018.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [6] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. In *NeurIPS*, 2018.
- [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [8] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Cite-seer, 2009.
- [9] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012.
- [10] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *ICLR*, 2017.
- [11] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [12] Wenbin Li, Lei Wang, Jinglin Xu, Jing Huo, Gao Yang, and Jiebo Luo. Revisiting local descriptor based image-to-class measure for few-shot learning. In *CVPR*, 2019.
- [13] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.

- [14] Takeru Miyato, Shin-ichi Maeda, Shin Ishii, and Masanori Koyama. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [15] Taesik Na, Jong Hwan Ko, and Saibal Mukhopadhyay. Cascade adversarial machine learning regularized with a unified embedding. In *ICLR*, 2018.
- [16] Ingram Olkin and Friedrich Pukelsheim. The distance between two random vectors with given dispersion matrices. *Linear Algebra and its Applications*, 48:257–263, 1982.
- [17] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. *ICLR*, 2017.
- [18] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *NeurIPS*, pages 5014–5026, 2018.
- [19] Chuanbiao Song, Kun He, Liwei Wang, and John E. Hopcroft. Improving the generalization of adversarial training with domain adaptation. In *ICLR*, 2019.
- [20] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *CVPR*, pages 403–412, 2019.
- [21] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *ICLR*, volume 1050, page 11, 2019.
- [22] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NeurIPS*, pages 3630–3638, 2016.
- [23] Kaidi Xu, Sijia Liu, Pu Zhao, Pin-Yu Chen, Huan Zhang, Deniz Erdogmus, Yanzhi Wang, and Xue Lin. Structured adversarial attack: Towards general implementation and better interpretability. In *ICLR*, 2019.
- [24] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. *CVPR*, 2018.
- [25] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan. Theoretically principled trade-off between robustness and accuracy. In *ICML*, 2019.
- [26] Stephan Zheng, Yang Song, Thomas Leung, and Ian J. Goodfellow. Improving the robustness of deep neural networks via stability training. In *CVPR*, pages 4480–4488, 2016.