
Robust Representation Learning via Asymmetric Negative Contrasting and Reverse Attention

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep neural networks are vulnerable to adversarial noise. Adversarial training (AT)
2 has been demonstrated to be the most effective defense strategy to protect neural
3 networks from being fooled. However, we find AT omits to learning robust features,
4 resulting in poor performance of adversarial robustness. To address this issue, we
5 highlight two characteristics of robust representation: (1) *exclusion: the feature of*
6 *natural examples keeps away from that of other classes;* (2) *alignment: the feature*
7 *of natural and corresponding adversarial examples is close to each other.* These
8 motivate us to propose a generic framework of AT to gain robust representation,
9 by the asymmetric negative contrast and reverse attention. Specifically, we design
10 an asymmetric negative contrast based on predicted probabilities and generate
11 adversarial negative examples by the targeted attack, to push away examples of
12 different classes in the feature space. Moreover, we propose to weight feature by
13 parameters of the linear classifier as the reverse attention, to obtain class-aware
14 feature and pull close the feature of the same class. Empirical evaluations on three
15 benchmark datasets show our methods greatly advance the robustness of AT and
16 achieve the state-of-the-art performance.

17 1 Introduction

18 Deep neural networks (DNNs) have achieved great success in academia and industry, but they
19 are easily fooled by carefully crafted adversarial examples to output incorrect results [13], which
20 leads to potential threats and insecurity in application. Given a well-trained DNN and a natural
21 example, an adversarial example can be generated by adding small perturbation that is invisible to
22 the human eyes to the natural example. The natural example can be correctly classified before the
23 perturbation and the adversarial example is incorrectly classified after the perturbation. In recent
24 years, there are many researches exploring the generation of adversarial examples to cheat models in
25 various fields, including image classification [13, 26, 5, 9], object detection [33, 8], natural language
26 processing [27, 2], semantic segmentation [28, 25], etc. The vulnerability of DNNs has aroused
27 common concerns on adversarial robustness.

28 Many empirical defense methods have been proposed to protect DNNs from adversarial perturbation,
29 such as adversarial training (AT) [26, 36, 30, 18, 39, 37, 31], image denoising [24], defensive
30 distillation [38, 6] and so on. The mainstream view is that AT is the most effective defense, which has
31 a training process of a two-sided game. The "attacker" crafts perturbation dynamically to generate
32 adversarial data to cheat the "defender", and the "defender" minimizes the loss function against
33 adversarial samples to improve robustness of models. Existing work [38, 6, 37, 11, 18, 20, 39] has
34 improved the effectiveness of AT in many aspects, but few studies pay attention to learning robust
35 feature. The overlook may lead to potential threats in the feature space of AT models, which does
36 harm to robust classification. Besides, there are no criteria for robust feature. In addition, adversarial

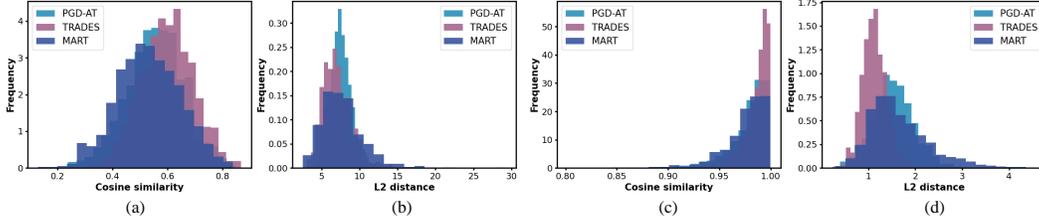


Figure 1: Frequency histograms of the L_2 distance and cosine similarity of the feature that belongs to natural examples, AEs and OEs. The four figures show the cosine similarity of the feature between natural examples and OEs (a), the L_2 distance of the feature between natural examples and OEs (b), the cosine similarity of the feature between natural examples and AEs (c), the L_2 distance of the feature between natural examples and AEs (d), respectively. The feature denotes the feature vector z before the linear layer. We train ResNet-18 [15] models on CIFAR-10 [22] with three AT methods: PGD-AT [26], TRADES [36] and MART [30]. In the calculation, we use all samples labeled as class 0 in the test set as natural examples and generate AEs by PGD-10 [26].

37 contrastive learning (ACL) and robust feature selection (RFS) are techniques to optimize feature
 38 distribution. ACL [21, 12, 35] is a kind of contrast learning (CL) [7, 17, 14] that extends to AT. RFS
 39 mostly modifies the architecture of models [32, 1, 34] to select important feature. However, the target
 40 problems of them are not to learn robust feature.

41 To demonstrate AT is indeed deficient in the representation which causes limited adversarial robust-
 42 ness, we conduct a simple experiment. We choose the L_2 distance and cosine similarity as metrics.
 43 And we measure the distance and similarity of the feature between natural examples, adversarial
 44 examples (AEs) and examples of other classes (OEs). The frequency histograms of the distance and
 45 similarity is shown in Figure 1. Figure 1 (a) and Figure 1 (b) show that the cosine similarity of the
 46 feature between natural examples and OEs shows a Gaussian distribution between 0.4 and 0.8, and
 47 the L_2 distance shows a skewed distribution between 2.0 and 12.0, which indicates there are very
 48 close pairs of natural examples and OEs that are not distinguished in the feature space. In Figure 1
 49 (c) and Figure 1 (d), it is shown that there are a skewed distribution between 0.9 and 0.99 for the
 50 cosine similarity of the feature between natural examples and AEs, and a skewed distribution between
 51 0.5 and 2.5 for the L_2 distance, which indicates that the feature of natural examples and AEs is not
 52 adequately aligned. Thus, there is still large room for optimization of the feature of AT.

53 Based on the observation, we propose two characteristics of robust feature: *exclusion: the feature*
 54 *of natural examples keeps away from that of other classes; alignment: the feature of natural and*
 55 *corresponding adversarial samples is close to each other. First, exclusion confirms the separability*
 56 *between different classes and avoids confusion in the feature space, which makes it hard to fool the*
 57 *model because the feature of different classes keep a large distance. Second, alignment insures the*
 58 *feature of natural examples is aligned with adversarial one, which guarantees the predicted results of*
 59 *the natural and adversarial examples of the same instances are also highly consistent. And it helps to*
 60 *narrow the gap between robust accuracy and clean accuracy.*

61 To address the issue, we propose an AT framework to concentrate on robust representation with the
 62 guidance of the two characteristics. Specifically, we suggest two strategies to meet the characteristics,
 63 respectively. Treat a natural example and corresponding AE as a positive pair (PP), and treat a natural
 64 example and corresponding OE as a negative pair (NP). For *exclusion*, we propose an asymmetric
 65 negative contrast based on predicted probabilities, which freezes natural examples and pushes away
 66 OEs by reducing the confidence of predicted class when predicted classes of NPs are consistent. In
 67 particular, we find OEs generated by the targeted attack are more beneficial for correct classification
 68 than those selected carefully. For *alignment*, we use the reverse attention to weight the feature of PPs
 69 by partial parameters of the linear classifier, which contains the importance of feature to target classes
 70 during classification. Because the feature of the same class gets the same weighting and feature of
 71 different classes is weighted disparately, PPs are aligned and each example of PPs becomes close
 72 to each other in the feature space. Empirical evaluations show that AT methods combined with our
 73 framework can greatly enhance robustness, which means the neglect of learning robust feature is
 74 one of the main reasons for poor robust performance of AT. In a word, we propose a generic AT
 75 framework with the Asymmetric Negative Contrast and Reverse Attention (ANCRA), to learn robust
 76 representation and advance robustness. Our main contributions are summarized as follows:

- 77 • We suggest improving adversarial training from the perspective of learning robust feature,
78 and two characteristics are highlighted as criteria of optimizing robust representation.
- 79 • We propose a generic framework of adversarial training, termed as ANCRA, to obtain robust
80 feature by the asymmetric negative contrast and reverse attention, with the guidance of two
81 characteristics of robust feature. It can be easily combined with other defense methods.
- 82 • Empirical evaluations show our framework can obtain robust feature and greatly improve
83 adversarial robustness, which achieves the of state-of-the-art performances on CIFAR-10,
84 CIFAR-100 and Tiny-ImageNet.

85 2 Related work

86 **Adversarial training** Madry et al. [26] propose PGD attack and PGD-based adversarial training,
87 forcing the model to correctly classify adversarial samples within the epsilon sphere during training
88 to obtain robustness, which is the pioneer of adversarial learning. Zhang et al. [36] propose to learn
89 both natural and adversarial samples and reduce the divergence of classification distribution of both
90 to reduce the difference between robust accuracy and natural accuracy. Wang et al. [30] find that
91 misclassified samples during training have a negative impact on robustness significantly, and propose
92 to improve the model’s attention to misclassification by adaptive weights. Zhang et al. [37] propose
93 to replace fixed attack steps with attack steps that just cross the decision boundary, and improved the
94 natural accuracy by appropriately reducing the number of attack iterations. Huang et al. [18] replace
95 labels with soft labels predicted by the model and adaptively reduce the weight of misclassification
96 loss to alleviate robust overfitting problem. Dong et al. [11] also propose a similar idea of softening
97 label and explain the different effects of hard and soft labels on robustness by investigating the
98 memory behavior of the model for random noisy labels. Chen et al. [6] propose random weight
99 smoothing and self-training based on knowledge distillation, which greatly improve the natural and
100 robust accuracy. Zhou et al. [39] embed a label transition matrix into models to infer natural labels
101 from adversarial noise. However, little work has been done to improve AT from the perspective
102 of robust feature learning. Our work shows AT indeed has defects in the feature distribution, and
103 strategies proposed to learn robust feature can greatly advance robustness, which indicates the neglect
104 of robust representation results in poor robust performance of AT.

105 **Adversarial contrastive learning** Kim et al. [21] propose an adversarial training method of
106 maximizing and minimizing the contrastive loss. Fan et al. [12] notice that the robustness of ACL
107 relies on fine-tuning, and pseudo labels and high-frequency information can advance robustness. Kucer
108 et al. [23] find that the direct combination of self-supervised learning and AT penalizes non-robust
109 accuracy. Bui et al. [3] propose some strategies to select positive and negative examples based on
110 predicted classes and labels. Yu et al. [35] find the instance-level identity confusion problem brought
111 by positive contrast and address it by asymmetric methods. The idea of these methods motivates us
112 to further consider how to obtain robust feature by contrast mechanism. We design a new negative
113 contrast to push away NPs and mitigate the confusion caused by negative contrast.

114 **Robust feature selection** Xiao et al. [32] take the maximum k feature values in each activation
115 layer to increase adversarial robustness. Zoran et al. [40] use a spatial attention mechanism to identify
116 important regions of the feature map. Bai et al. [1] propose to suppress redundant feature channels and
117 dynamically activate feature channels with the parameters of additional components. Yan et al. [34]
118 propose to amplify the top-k activated feature channels. Existing work has shown enlarging import
119 feature channels is beneficial for robustness, but most approaches rely on extra model components and
120 do not explain the reason. We proposes the reverse attention to weight feature by class information
121 without any extra components, and explain it by *alignment* of feature.

122 3 Methodology

123 This section explains the instantiation of the our AT framework from the perspective of the two
124 characteristics of robust feature. To meet *exclusion*, we design an asymmetric negative contrast based
125 on predicted probabilities and propose to craft OEs by the targeted attack, to push away the feature of
126 NPs. To confirm *alignment*, we propose the reverse attention to weight the feature of the same class,
127 by the corresponding weight of target class in parameters of the linear classifier, so that the feature of
128 PPs is aligned and the gap of the feature between natural examples and AEs becomes small.

129 **3.1 Notations**

130 In this paper, capital letters indicate random variables or vectors, while lowercase letters represent
 131 their realisations. We define the function for classification as $f(\cdot)$. It can be parameterized by
 132 DNNs. $Linear(\cdot)$ is the linear classifier with a weight of Ω (\mathbb{C} , \mathbb{R}), in which \mathbb{C} denotes the class
 133 number and \mathbb{R} denotes the channel number of the feature map. $g(\cdot)$ is the feature extractor, i.e.,
 134 the rest model without $Linear(\cdot)$. Let $\mathcal{B} = \{x_i, y_i\}_i^N$ be a batch of natural samples where x_i
 135 is labeled by y_i . Given an adversarial transformation \mathcal{T}_a from an adversary \mathcal{A} (e.g., PGD attack
 136 in [26]), and a strategy \mathcal{T}_o for selection or generation of OEs. For data, we consider a positive pair
 137 $PP = \{x_i, x_i^a | x_i \in \mathcal{B}, x_i^a = \mathcal{T}_a(x_i)\}_i^N$, and a negative pair $NP = \{x_i, x_i^o | x_i \in \mathcal{B}, x_i^o = \mathcal{T}_o(x_i)\}_i^N$. Let
 138 $\mathbb{N}(x, \epsilon)$ represent the neighborhood of $x : \{\tilde{x} : \|\tilde{x} - x\| \leq \epsilon\}$, where ϵ is the perturbation budget. For
 139 an input x_i , we consider its feature z_i before $Linear(\cdot)$, the probability vector $p_i = softmax(f(x_i))$
 140 and predicted class $h_i = argmax(p_i)$, respectively.

141 **3.2 Adversarial training with asymmetric negative contrast**

142 Firstly, we promote AT to learn robust representation that meets *exclusion*. We notice that ACL has
 143 the contrastive loss [29] to maximize the consistency between PPs and to minimize the consistency
 144 between NPs. Motivated by the contrast mechanism, we consider to design a new negative-contrast
 145 term and combine it with AT loss, which creates a repulsive action between NPs when minimize the
 146 whole loss. Thus, we propose a generic pattern of AT loss with a negative contrast. Let TRADES
 147 [36] represent AT in the following paper as a example.

$$\mathcal{L}^{CAL}(x, y, x^a, x^o) = \mathcal{L}^{TRADES} + Sim(x, x^o) = \mathcal{L}_{CE}(x, y) + \mathcal{D}_{KL}(x, x^a) + Sim(x, x^o), \quad (1)$$

148 Where x denotes natural examples with labels y , x^a are AEs generated by untargeted PGD [26],
 149 x^o are negative examples of other classes (OEs), Sim is a similarity function, \mathcal{L}_{CE} denotes the
 150 cross-entropy loss and \mathcal{D}_{KL} denotes divergence of Kullback-Leibler. AEs generated by maximizing
 151 \mathcal{L}_{CE} typically have wrong predicted classes, given by:

$$x_{t+1}^a := \mathbf{\Pi}_{\mathbb{N}(x, \epsilon)}(x_t^a + \epsilon \text{sign}(\nabla_x \mathcal{L}_{CE}((f(x_t^a), y))), \quad (2)$$

152 where ϵ denotes the L_∞ -norm of perturbation, x_t^a denotes adversarial positive samples after the t th
 153 attack iteration, $\mathbf{\Pi}$ denotes a clamp function, $sign$ denotes a sign function and $\nabla_x \mathcal{L}_{CE}$ denotes the
 154 gradient of \mathcal{L}_{CE} with respect to x . When minimizing the loss in Eq 1, \mathcal{L}^{TRADES} learns to classify
 155 natural examples and AEs correctly, and additional negative contrast prompts the inconsistency of
 156 NPs, which keeps the feature of NPs far away from each other. The whole loss guides the model to
 157 learn correct classification from TRADES and push away NPs from each other to ensure *exclusion*.
 158 Although we have a generic pattern of AT loss with a negative contrast, there are several problems
 159 about details to address. To refine the negative contrast and address problems, we further propose a
 160 method to calculate the negative contrast and strategy to generate OEs.

161 **3.2.1 Asymmetric negative contrast based on probabilities**

162 The work in [35] has indicated that when the predicted classes of the adversarial positive examples
 163 (i.e., AEs) and negative samples (i.e., OEs) are the same, the positive contrast may lead to a conflict
 164 between the positive and negative contrast, resulting in wrong classification. On the basis, we find
 165 a similar conflict can also be caused by the negative contrast when the predicted classes of AEs
 166 and OEs are different, which we named by class confusion. As shown in Figure 2, when AEs and
 167 OEs have different predicted classes, natural examples are subject to the attraction of AEs and the
 168 repulsion of OEs at the same time. And it is likely to move near the decision boundary or even into
 169 the wrong class space under the actions, which does harm to *exclusion*.

170 In order to alleviate the problem of class confusion, We should reasonably control the effect of the
 171 repulsion of negative contrast between natural examples and OEs. we propose an asymmetric method
 172 of the negative contrast, $Sim^\alpha(x, x^o)$, to decouple the repulsive force into an one-side push from the
 173 natural example to the OE and an one-side push from the OE to the natural example, given by:

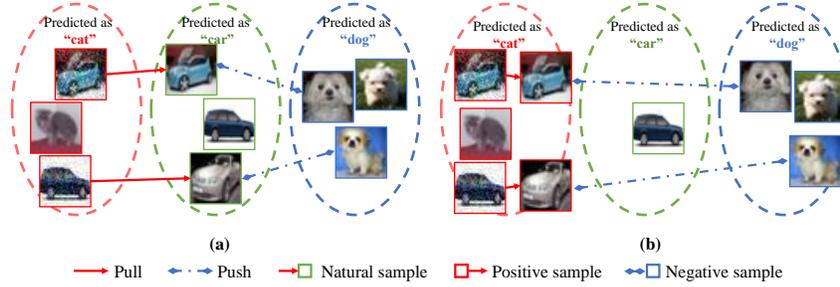


Figure 2: Illustrations of class confusion when the classes of positive examples (i.e., AEs) and negative examples (i.e., OEs) are different. (a) shows the normal situation before the optimization. (b) shows the situation of class confusion after the optimization. In each circle, data points have the same predicted class. In (a), AEs locate in the wrong predicted class different from natural example and OEs. The TRADES loss narrow the gap of classification between natural examples and AEs, and thus AEs in the wrong class pull natural examples to move toward the wrong class and the negative contrast pushes natural examples to leave from the original class. With these actions, natural examples come to the decision boundary and even into the wrong class easily as (b) shows.

$$\text{Sim}^\alpha(x, x^o) = \alpha \cdot \overline{\text{Sim}}(x, x^o) + (1 - \alpha) \cdot \overline{\text{Sim}}(x^o, x), \quad (3)$$

174 where $\overline{\text{Sim}}(x, x^o)$ denotes the one-sided similarity of x and x^o . When minimizing $\overline{\text{Sim}}(x, x^o)$, we
 175 stop the back-propagation gradient of x and only move x^o away from x . α denotes the weighting
 176 factor to adjust the magnitude of the two repulsive forces. When $\alpha = 0$, OEs are frozen and only the
 177 feature of natural samples is optimized to push far away from the feature of OEs. As α increases, the
 178 natural sample becomes more repulsive to the OE and the OE pushes the natural example less. To
 179 mitigate the class confusion problem, we should choose α that tends to 1 to reduce the repulsive force
 180 from the OE to the natural example, to prevent the natural example from being pushed into the wrong
 181 class. Experiments show that $\alpha = 1$ leads to the best performance provided in our supplementary
 182 material), which pushes away NPs by only pushing off OEs and follows what we have expected.

183 Then we propose the negative contrast based on predicted probabilities, $\text{Sim}_{cc}^\alpha(x, x^o)$, to measure
 184 the repulsive force of NPs pushing away from each other. It pushes away NPs by decreasing the
 185 corresponding probabilities of the predicted classes when the predicted classes of NPs are consistent.

$$\text{Sim}_{cc}^\alpha(x, x^o) = \frac{1}{\|\mathcal{B}_i\|} \sum_{i=1}^n \mathbb{I}(h_i = h_i^o) \cdot \left[\alpha \sqrt{\hat{p}_i(h_i) \cdot p_i^o(h_i)} + (1 - \alpha) \sqrt{p_i(h_i) \cdot \hat{p}_i^o(h_i)} \right], \quad (4)$$

186 where $\|\mathcal{B}_i\|$ denotes the batch size, $\mathbb{I}(\cdot)$ denotes the Indicator function and \hat{p} denotes freezing the
 187 back-propagation gradient of p . h_i and h_i^o denote the predicted classes of the NP. And p_i and p_i^o
 188 denote the probability vectors of the NP. Under the negative contrast, the model pushes the natural
 189 example in the direction away from the predicted class of the OE and push the OE in the direction
 190 away from the predicted class of the natural example when and only when two predicted classes of
 191 the NP are consistent. This ensures that the action of *exclusion* not only pushes away the feature of
 192 NPs in the feature space, but also reduces the probabilities of NPs in the incorrect class. Since the
 193 negative contrast has only directions to reduce the confidence and no explicit directions to increase
 194 the confidence, it does not create any actions to push the natural example into the feature space of
 195 wrong classes even in the scenario of class confusion, which can effectively alleviate the problem.

196 3.2.2 Generate negative samples by targeted attack

197 To obtain OEs, previous negative sampling strategies [19] simply screen natural samples and pick up
 198 the negatives from them, but rarely consider generating special negative samples to assist learning.
 199 We innovatively propose a strategy to craft OEs by the targeted attack: natural negative examples with
 200 labels that is different from those of natural examples are attacked to the labeled classes of natural
 201 examples by targeted PGD-10 [26], to manufacture hard negatives containing adversarial noise.

$$x_{t+1}^o := \prod_{\mathbb{N}(x^o, \epsilon)} (x_t^o - \epsilon \text{sign}(\nabla_{x^o} \mathcal{L}_{CE}((f(x_t^o), y))), \quad (5)$$

202 Where $\nabla_{x^o} \mathcal{L}_{CE}$ denotes the gradient of \mathcal{L}_{CE} with respect to x^o . By this strategy, clean OEs randomly
 203 chosen from other classes are attacked to the labeled classes of natural examples and become negative
 204 adversarial examples. The motivation makes intuitive sense. 1) The negative adversarial sample
 205 generated by the targeted attack will be classified as the labeled class of the natural example with
 206 high confidence, but its ground truth label is not that, which makes it a very hard negative sample
 207 and is beneficial for the negative contrast. 2) The negative adversarial sample contains adversarial
 208 noise, which is special feature that natural negative samples do not have. And this feature helps the
 209 model learn the paradigm of adversarial noise and improve the robust performance. In particular, we
 210 demonstrate that negative samples with adversarial noise do improve robustness better in Table 4.

211 3.3 Adversarial training with reverse Attention

212 Secondly, we continue to improve TRADES to learn robust representation that meets *alignment*.
 213 Consider the calculating process of the model $f(\cdot)$. First, the feature vector z is obtained by $g(x)$,
 214 and then the output vector Ωz is obtained by a linear mapping $Linear(z)$. Each element z_i in z
 215 represents the activation level of the feature channel that may be helpful for classification, with larger
 216 values representing more feature information extracted from that channel; $\omega^{i,j}$ in Ω represents the
 217 importance of the i th feature channel to the j th class, with higher values representing the greater
 218 contribution of the feature channel to the class. Motivated by [1, 34], we exploit the importance
 219 of feature channels to target classes to align the feature of examples of the same classes and pull
 220 close the feature of PPs, which is named by reverse attention. To be specific, we take the Hadamard
 221 product (Kronecker product) of partial weight of the classifier Ω^j and the feature vector z . It can
 222 weight feature channel by channel according to its contribution to being classified as the target class
 223 j , and gain a class-aware feature vector z' containing the information of the target class j .

$$z'_i = \begin{cases} z_i \odot \omega^{i,y}, & \text{(training phase)} \\ z_i \odot \omega^{i,h(x)}, & \text{(testing phase)} \end{cases} \quad (6)$$

224 where \odot denotes the Hadamard product operation, which is the method of multiplying two matrices of
 225 the same size element by element to obtain a new matrix of the same size. To ensure parameters used
 226 for weighting have the correct feature-to-class importance, we use the unweighted feature vector z to
 227 go through $Linear(\cdot)$ to obtain the auxiliary probability vector p , and z' to get the final probability
 228 vector p' . Finally, we use both p and p' to train the model. During the training phase, we use the true
 229 label y as an indicator to determine the importance of channels, i.e., $\Omega^j = \Omega^y$. And in the testing
 230 phase, since the true label is not available, we simply choose a sub-vector of the linear weight by the
 231 predicted class $h(x)$ as the importance of channels. We add the reverse attention to the last feature
 232 layer in the model, which generally contains two blocks. The model with the reverse attention does
 233 not need any extra modules, but module interactions are changed.

234 Let's make a detailed analysis and explanation of the principle of this method. The class information
 235 from labels guides the input image to be mapped from the feature to the classification vector during
 236 training, establishing an feature-to-class mapping relationship. In the model, the feature extractor
 237 captures the representation that is helpful for classification until the feature vector contains enough
 238 information that allows the classifier to classify the sample as the target class. Among all the modules,
 239 the classifier is the closest to labels and learns which feature channel plays an important role in being
 240 classified as the target class (i.e., the feature importance). Since the classifier is unique, the importance
 241 of the feature channels of one example is exactly the same with that of the other samples in the same
 242 class, benefiting the generalization and robustness of the model in the target class. We propose the
 243 reverse attention to utilize this information to improve feature rather than classification. The feature
 244 vectors are weighted by partial parameters of the linear layer that belong to the target class, which
 245 can change the activation of each channel adaptively according to the feature importance, acting as an
 246 attention with the guidance of the class information. After the attention, the important channels in the
 247 feature vector are boosted and the redundant channels are weakened, i.e., the information contributes
 248 to the target class will become larger and more significant, which is helpful for correct classification.
 249 Considering from the perspective of the feature distribution, the weighted feature has gained extra
 250 class information, which induces changes in the feature distribution. Feature vectors with the same
 251 target class get the same weighting, and thus the weighted feature becomes more similar. Moreover,
 252 feature vectors with different target classes are weighted according to different weights, and the
 253 weighted feature distributions become more inconsistent. Therefore, the reverse attention guides

254 the *alignment* of the feature of the examples in the same class, pulling the feature of PPs closer and
 255 pushing the feature of NPs far away, which benefits *alignment* and drops by to promote *exclusion* and
 256 classification. Aligned feature has similar activations in every feature channel, which helps the model
 257 narrows the gap between feature of natural examples and AEs.

258 4 Experiments

259 To demonstrate the effectiveness of the proposed approach, we show feature distribution of trained
 260 models firstly. Then we evaluate our framework against white-box attacks and adaptive attacks, and
 261 make a compare with other defense methods. We conduct experiments across different datasets
 262 and models. Because our methods are compatible with existing AT techniques and can be easily
 263 incorporated in a plug-and-play manner, we choose three baselines [26, 36, 30] to combine with our
 264 framework for evaluation: PGD-AT-ANCRA, TRADES-ANCRA, and MART-ANCRA.

265 4.1 Settings

266 **Implementation** On CIFAR-10 and CIFAR-100 [22], we train ResNet18 [15] with a weight
 267 decay of 2.0×10^{-4} . On Tiny-ImageNet [10], we use PreActResNet18 [16] with a weight decay of
 268 5.0×10^{-4} . We adopt the SGD optimizer with a learning rate of 0.01, a momentum of 0.9, epochs of
 269 120 and a batch size of 128 as [30]. For the trade-off hyperparameters β , we use 6.0 in TRADES
 270 and 5.0 in MART, following the original setting in their papers. For other hyperparameters, we tune
 271 the values based on TRADES-ANCRA. We generate adversarial example for training by L_∞ -norm
 272 PGD [26], with a step size of 0.007, an attack iterations of 10 and perturbation budget of 8/255. We
 273 use single NVIDIA A100 and two GTX 2080 Ti in the experiments.

274 **Baseline** We compare the proposed PGD-AT-ANCRA, TRADES-ANCRA, and MART-ANCRA
 275 with the popular baselines: PGD-AT [26], TRADES [36], MART [30] and SAT [18]. Moreover, we
 276 also choose three state-of-the-art methods: AWP [31], S2O [20] and UDR [4]. We keep the same
 277 settings among all the baselines with our settings and follow their original hyperparameters.

278 **Evaluation** We choose several adversarial attacks to attack the target models, including PGD [26],
 279 FGSM [13], C&W [5] and AutoAttack [9] which is a powerful and reliable attack and an ensemble
 280 attack with three white-box attacks and one black-box attack. We notice that our methods use the
 281 auxiliary probability vector p in the training and testing phase, so we design two scenarios: 1) train
 282 with p and test without p ; 2) train with p and test with p . 1) denotes evaluation against white-box
 283 attacks and 2) denotes evaluation against adaptive attacks. Following the default setting of AT, the
 284 max perturbation strength is set as 8. / 255. for all attack methods under the L_∞ . The attack iterations
 285 of PGD and C&W is 40 (i.e., PGD-40), and the step size of FGSM is 8. / 255. unlike 0.007 for other
 286 attacks. The clean accuracy and robust accuracy are used as the evaluation metrics.

287 4.2 Comparison results of feature distribution

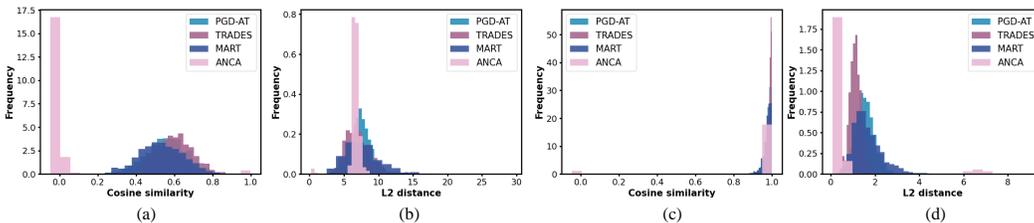


Figure 3: Frequency histograms of the L_2 distance and cosine similarity of feature of natural examples, AEs and OEs. We train ResNet-18 models on CIFAR-10 with four defense techniques: PDG-AT, TRADES, MART and TRADES-ANCRA. Other details are the same with Figure 1

288 Frequency histograms of feature distribution is shown in Figure 3. It is shown that our methods can
 289 greatly improve feature distribution, which follows the characteristics of *exclusion* and *alignment*. In
 290 Figure 3 (a) and Figure 3 (b), it shows that the cosine similarity of the model trained by our method
 291 between natural examples and OEs shows a skewed distribution between -0.05 and 0.1, and the L_2

292 distance with our method shows a Gaussian distribution between 5.5 and 10.0, which indicates natural
 293 examples and OEs have been fully distinguished in the feature space and *exclusion* has been met. In
 294 Figure 3 (c) and Figure 3 (d), it shows that in the model trained by our method there are a uniform
 295 distribution between 0.95 and 0.99 for the cosine similarity of the feature between natural examples
 296 and AEs, and a skewed distribution between 0.05 and 1.5 for the L_2 distance of the feature, which
 297 indicates the feature between natural examples and AEs is very close to each other and *alignment* has
 298 been confirmed. Thus, our framework successfully helps AT to obtain robust feature.

299 4.3 Comparison results against white-box attacks

Table 1: Robustness (%) against white-box attacks. Nat denotes clean accuracy. PGD denotes robust accuracy against PGD-40. FGSM denotes robust accuracy against FGSM. C&W denotes robust accuracy against C&W. AA denotes robust accuracy against AutoAttack. Mean denotes average robust accuracy against these four attacks. We show the most successful defense with **bold**.

Defense	CIFAR-10						CIFAR-100					
	Nat	PGD	FGSM	C&W	AA	Mean	Nat	PGD	FGSM	C&W	AA	Mean
PGD-AT	80.90	44.35	58.41	46.72	42.14	47.91	56.21	19.41	30.00	41.76	17.76	27.23
TRADES	78.92	48.40	59.60	47.59	45.44	50.26	53.46	25.37	32.97	43.59	21.35	30.82
MART	79.03	48.90	60.86	45.92	43.88	49.89	53.26	25.06	33.35	38.07	21.04	29.38
SAT	63.28	43.57	50.13	47.47	39.72	45.22	42.55	23.30	28.36	41.03	18.73	27.86
AWP	76.38	48.88	57.47	48.22	44.65	49.81	54.53	27.35	34.47	44.91	21.98	31.18
S2O	40.09	24.05	29.76	47.00	44.00	36.20	26.66	13.11	16.83	43.00	21.00	23.49
UDR	57.80	39.79	45.02	46.92	34.73	41.62	33.63	20.61	24.19	33.77	16.41	23.75
PGD-AT-ANCRA	85.10	89.03	87.00	89.23	59.15	81.10	59.73	58.10	58.45	58.58	34.44	52.39
TRADES-ANCRA	81.70	82.96	82.74	83.01	59.70	77.10	53.73	51.24	52.17	52.55	35.81	47.94
MART-ANCRA	84.88	88.56	87.95	88.77	59.62	81.23	60.10	58.40	58.74	59.41	35.05	52.90

300 We train ResNet-18 by different defense on CIFAR-10 and CIFAR-100 to evaluate them under
 301 white-box attacks. And more results in PreActResNet18 on Tiny-ImageNet are provided in our
 302 supplementary material. The results on CIFAR-10 and CIFAR-100 are shown in Table 1. First, on
 303 CIFAR-10, our approaches improve the clean accuracy of based approaches by 5.2%, 3.2% and 5.9%,
 304 and also improves the robust performance under all the attacks (e.g., increase by 44.7%, 34.6% and
 305 39.7% against PGD). Compared with state-of-the-art defense, the robust accuracy against different
 306 attacks of our methods is almost two times as large than theirs (e.g., 81.23% VS 49.81%). Second, on
 307 CIFAR-100, our approaches also greatly improve the robustness and advance the clean accuracy. The
 308 clean accuracy of our methods has been increased by 3.5%, 0.3% and 6.8% compared with based
 309 methods, and the lowest average robust accuracy of ours is larger than the best one among other
 310 methods by 16.8%. In general, our three approaches gain the best performance both in the natural and
 311 attacked scenarios. To our surprise, MART-ANCRA and PGD-ANCRA rather than TRADES-ANCRA
 312 gain the best performance in a lot of cases without hyper-parameter tuning. Besides, our approaches
 313 not only improves robustness but also enhances clean accuracy, though there is always a trade-off
 314 between clean and robust accuracy. These results indicate that our approaches can vastly boost the
 315 robustness of models against white-box attacks.

316 4.4 Comparison results against adaptive attacks

317 We train several ResNet18 models on CIFAR-10 by PGD-AT-ANCRA, TRADES-ANCRA, MART-
 318 ANCRA and test the same models without p . In addition, we report vanilla based approaches as
 319 baseline. Results are in Table 2. It indicates that our approaches can still maintain superb performance
 320 after adaptive attacks, e.g., the robust accuracy against PGD of our methods without p are larger than
 321 those of baseline by 13.28%, 10.08% and 8.06%.

322 4.5 Ablation studies

323 **Two defense methods.** We train four models by TRADES, TRADES with the asymmetric negative
 324 contrast (TRADES-ANC), TRADES with the reverse attention (TRADES-RA) and TRADES-
 325 ANCRA, respectively. The results of evaluation against adaptive attacks are shown in Table 3. First,

Table 2: Robustness(%) of ResNet-18 trained with our approaches and attacked with or without p .

Approach	Nat	Attack with p			Attack without p		
		PGD	FGSM	C&W	PGD	FGSM	C&W
Vanilla TRADES	78.92	\	\	\	48.40	59.60	47.59
TRADES-ANCRA	81.70	61.68	61.56	72.36	82.96	82.74	83.01
Vanilla PGD-AT	80.90	\	\	\	44.35	58.41	46.72
PGD-AT-ANCRA	85.10	54.43	58.23	66.36	89.03	87.00	89.23
Vanilla MART	79.09	\	\	\	48.90	60.86	45.92
MART-ANCRA	84.88	56.96	60.43	71.06	88.56	87.95	88.77

326 when incorporating the asymmetric negative contrast only, the performance of robustness against all
 327 the attacks and clean accuracy have been improved compared with vanilla TRADES (e.g., 48.36% VS
 328 54.18% against PGD-40). Next, when incorporating the reverse attention only, the performance on
 329 clean and adversarial data is also improved greatly compared with TRADES (e.g., 48.36% VS 61.69%
 330 against PGD-40). Thus, it shows each method contributes to robustness and generalization. Besides,
 331 when TRADES-ANCRA is compared with TRADES-RA, the clean accuracy and robust accuracy
 332 against all the attacks except AA have been enhanced, which indicates that the two strategies are
 333 compatible and the combination can alleviate the side effect of independent methods.

334 **Strategy of negative samples** We compare our strategy of the targeted attack with other strategies
 335 to select negative samples, including Random, Soft-LS and Hard-LS proposed by Bui et al. [3]. The
 336 details of them are provided in our supplementary material. The results are shown in Table 4. To make
 337 a comprehensive compare, we show results of both the best models and last models with different
 338 strategies. It shows that our strategy have the best performance of robustness and clean accuracy in
 339 the last models, and achieve the best robust accuracy in the best models.

Table 3: Clean and robust accuracy (%) of ResNet-18 trained by TRADES, TRADES-ANC, TRADES-RA and TRADES-ANCRA on CIFAR-10 against various attacks.

Defense	Nat	PGD	FGSM	C&W	AA
TRADES	78.92	48.40	59.60	47.59	45.44
TRADES-ANC	80.77	54.18	63.44	49.84	48.51
TRADES-RA	80.46	61.59	61.48	72.15	61.02
TRADES-ANCRA	81.70	61.68	61.56	72.36	59.70

Table 4: Results of the best and last with four strategies of negative example. Best- denotes results in the best models and Last- denotes results in the last models. We show the best results with **bold**.

Strategy	Best-Nat	Best-PGD	Last-Nat	Last-PGD
Random	81.44	62.64	81.78	61.71
Soft-LS	82.10	61.83	80.62	58.47
Hard-LS	82.30	62.53	82.13	60.98
Targeted attack	81.36	63.08	82.18	62.02

340 5 Conclusion

341 This work addresses an overlook of robust representation learning in the adversarial training by a
 342 generic AT framework with the asymmetric negative contrast and reverse attention. We propose
 343 two characteristics of robust feature to guide the improvement of AT, i.e., *exclusion* and *alignment*.
 344 Specifically, the asymmetric negative contrast based on probabilities fixes natural examples, and only
 345 pushes away adversarial examples of other classes in the feature space. Besides, the reverse attention
 346 weights feature by parameters of the linear classifier, to provide class information and align feature of
 347 the same class. Our framework can be used in a plug-and-play manner with other defense methods.
 348 Analysis and empirical evaluations demonstrate that our framework can obtain robust feature and
 349 greatly improve robustness and generalization.

References

- 350
- 351 [1] Yang Bai, Yuyuan Zeng, Yong Jiang, Shu-Tao Xia, Xingjun Ma, and Yisen Wang. Improving adversarial
352 robustness via channel-wise activation suppressing. *arXiv preprint arXiv:2103.08307*, 2021.
- 353 [2] Nicholas Boucher, Iliia Shumailov, Ross Anderson, and Nicolas Papernot. Bad characters: Imperceptible
354 nlp attacks. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1987–2004, 2022. doi:
355 10.1109/SP46214.2022.9833641.
- 356 [3] Anh Bui, Trung Le, He Zhao, Paul Montague, Seyit Camtepe, and Dinh Phung. Understanding and achiev-
357 ing efficient robustness with adversarial supervised contrastive learning. *arXiv preprint arXiv:2101.10027*,
358 2021.
- 359 [4] Tuan Anh Bui, Trung Le, Quan Hung Tran, He Zhao, and Dinh Q. Phung. A unified wasserstein
360 distributional robustness framework for adversarial training. *CoRR*, abs/2202.13437, 2022. URL <https://arxiv.org/abs/2202.13437>.
- 362 [5] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *2017 IEEE
363 Symposium on Security and Privacy (SP)*, pages 39–57, 2016.
- 364 [6] Tianlong Chen, Zhenyu Zhang, Sijia Liu, Shiyu Chang, and Zhangyang Wang. Robust overfitting may be
365 mitigated by properly learned smoothing. In *International Conference on Learning Representations*,
366 2021.
- 367 [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
368 contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of
369 the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning
370 Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/v119/
371 chen20j.html](https://proceedings.mlr.press/v119/chen20j.html).
- 372 [8] Xiangning Chen, Cihang Xie, Mingxing Tan, Li Zhang, Cho-Jui Hsieh, and Boqing Gong. Robust and
373 accurate object detection via adversarial learning. In *Proceedings of the IEEE/CVF conference on computer
374 vision and pattern recognition*, pages 16622–16631, 2021.
- 375 [9] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
376 diverse parameter-free attacks. In *International Conference on Machine Learning*, 2020.
- 377 [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
378 image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255,
379 2009.
- 380 [11] Yinpeng Dong, Ke Xu, Xiao Yang, Tianyu Pang, Zhijie Deng, Hang Su, and Jun Zhu. Exploring
381 memorization in adversarial training. *arXiv preprint arXiv:2106.01606*, 2021.
- 382 [12] Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. When does contrastive learning
383 preserve adversarial robustness from pretraining to finetuning? *Advances in neural information processing
384 systems*, 34:21480–21492, 2021.
- 385 [13] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
386 examples. *CoRR*, abs/1412.6572, 2014.
- 387 [14] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya,
388 Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray
389 Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent a new approach to self-supervised
390 learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*,
391 NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 392 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
393 In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
394 doi: 10.1109/CVPR.2016.90.
- 395 [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks.
396 In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages
397 630–645, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46493-0.
- 398 [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised
399 visual representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
400 (CVPR)*, pages 9726–9735, 2020. doi: 10.1109/CVPR42600.2020.00975.

- 401 [18] Lang Huang, Chao Zhang, and Hongyang Zhang. Self-adaptive training: beyond empirical risk minimiza-
402 tion. *Advances in neural information processing systems*, 33:19365–19376, 2020.
- 403 [19] Tri Huynh, Simon Kornblith, Matthew R Walter, Michael Maire, and Maryam Khademi. Boosting
404 contrastive self-supervised learning with false negative cancellation. In *Proceedings of the IEEE/CVF*
405 *winter conference on applications of computer vision*, pages 2785–2795, 2022.
- 406 [20] Gaojie Jin, Xinping Yi, Wei Huang, Sven Schewe, and Xiaowei Huang. Enhancing adversarial training
407 with second-order statistics of weights. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
408 *and Pattern Recognition*, pages 15273–15283, 2022.
- 409 [21] Minseon Kim, Jihoon Tack, and Sung Ju Hwang. Adversarial self-supervised contrastive learning. *Advances*
410 *in Neural Information Processing Systems*, 33:2983–2994, 2020.
- 411 [22] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 412 [23] Michal Kucer, Diane Oyen, and Garrett Kenyon. When does visual self-supervision aid adversarial training
413 in improving adversarial robustness?
- 414 [24] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against
415 adversarial attacks using high-level representation guided denoiser. In *2018 IEEE/CVF Conference on*
416 *Computer Vision and Pattern Recognition*, pages 1778–1787, 2018. doi: 10.1109/CVPR.2018.00191.
- 417 [25] Yawei Luo, Ping Liu, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Category-level adversarial
418 adaptation for semantic segmentation using purified features. *IEEE Transactions on Pattern Analysis and*
419 *Machine Intelligence*, 44(8):3940–3956, 2022. doi: 10.1109/TPAMI.2021.3064379.
- 420 [26] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards
421 deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- 422 [27] John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. TextAttack: A framework
423 for adversarial attacks, data augmentation, and adversarial training in NLP. In *Proceedings of the 2020*
424 *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages
425 119–126, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.
426 emnlp-demos.16. URL <https://aclanthology.org/2020.emnlp-demos.16>.
- 427 [28] Federico Nesti, Giulio Rossolini, Saasha Nair, Alessandro Biondi, and Giorgio Buttazzo. Evaluating the
428 robustness of semantic segmentation for autonomous driving against real-world adversarial patch attacks.
429 In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2280–2289,
430 2022.
- 431 [29] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
432 coding. *ArXiv*, abs/1807.03748, 2018.
- 433 [30] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adver-
434 sarial robustness requires revisiting misclassified examples. In *International Conference on Learning*
435 *Representations*, 2020.
- 436 [31] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization.
437 *Advances in Neural Information Processing Systems*, 33:2958–2969, 2020.
- 438 [32] Chang Xiao, Peilin Zhong, and Changxi Zheng. Enhancing adversarial defense by k-winners-take-all.
439 *arXiv preprint arXiv:1905.10510*, 2019.
- 440 [33] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples
441 for semantic segmentation and object detection. In *Proceedings of the IEEE international conference on*
442 *computer vision*, pages 1369–1378, 2017.
- 443 [34] Hanshu Yan, Jingfeng Zhang, Gang Niu, Jiashi Feng, Vincent Tan, and Masashi Sugiyama. Cifs: Improving
444 adversarial robustness of cnns via channel-wise importance-based feature selection. In *International*
445 *Conference on Machine Learning*, pages 11693–11703. PMLR, 2021.
- 446 [35] Qiyang Yu, Jieming Lou, Xianyuan Zhan, Qizhang Li, Wangmeng Zuo, Yang Liu, and Jingjing Liu.
447 Adversarial contrastive learning via asymmetric infonce. In *Computer Vision—ECCV 2022: 17th European*
448 *Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part V*, pages 53–69. Springer, 2022.
- 449 [36] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theo-
450 retically principled trade-off between robustness and accuracy. In *International conference on machine*
451 *learning*, pages 7472–7482. PMLR, 2019.

- 452 [37] Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankanhalli.
453 Attacks which do not kill training make adversarial learning stronger. In *International conference on*
454 *machine learning*, pages 11278–11287. PMLR, 2020.
- 455 [38] Shiji Zhao, Jie Yu, Zhenlong Sun, Bo Zhang, and Xingxing Wei. Enhanced accuracy and robustness via
456 multi-teacher adversarial distillation. In *European Conference on Computer Vision*, 2022.
- 457 [39] Dawei Zhou, Nannan Wang, Bo Han, and Tongliang Liu. Modeling adversarial noise for adversarial
458 training. In *International Conference on Machine Learning*, pages 27353–27366. PMLR, 2022.
- 459 [40] Daniel Zoran, Mike Chrzanowski, Po-Sen Huang, Sven Gowal, Alex Mott, and Pushmeet Kohli. Towards
460 robust image classification using sequential attention models. In *Proceedings of the IEEE/CVF conference*
461 *on computer vision and pattern recognition*, pages 9483–9492, 2020.