

CONFLICT-AWARE ADVERSARIAL TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

Adversarial training is the most effective method to obtain adversarial robustness for deep neural networks by directly involving adversarial samples in the training procedure. To obtain an accurate and robust model, the weighted-average method is applied to optimize standard loss and adversarial loss simultaneously. In this paper, we argue that the weighted-average method does not provide the best tradeoff for the standard performance and adversarial robustness. We argue that the failure of the weighted-average method is due to the conflict between the gradients derived from standard and adversarial loss, and further demonstrate such a conflict increases with attack budget theoretically and practically. To alleviate this problem, we propose a new trade-off paradigm for adversarial training with a conflict-aware factor for the convex combination of standard and adversarial loss, named **Conflict-Aware Adversarial Training (CA-AT)**. Comprehensive experimental results show that CA-AT consistently offers a superior trade-off between standard performance and adversarial robustness under the settings of adversarial training from scratch and parameter-efficient finetuning.

1 INTRODUCTION

Deep learning models have achieved exemplary performance across diverse application domains (He et al., 2017; Vaswani et al., 2017; Ouyang et al., 2022; Rombach et al., 2022; Radford et al., 2021). However, they remain vulnerable to adversarial samples produced by adversarial attacks (Goodfellow et al., 2014; Liu et al., 2016; Moosavi-Dezfooli et al., 2016). Deep learning models can easily be fooled into making mistakes by adding an imperceptible noise produced by adversarial attacks to the standard sample. To solve this problem, many methods have been proposed to improve the robustness against adversarial samples (Cai et al., 2018; Chakraborty et al., 2018; Madry et al., 2018), among which **adversarial training (AT)** has been proven to be the most effective strategy (Madry et al., 2018; Athalye et al., 2018; Qian et al., 2022; Bai et al., 2021). Specifically, AT aims to enhance model robustness by directly involving adversarial samples during training. They used adversarial examples to construct the adversarial loss functions for parameter optimization. The adversarial loss can be formulated as a min-max optimization objective, where the adversarial samples are generated by the inner maximization, and the model parameters are optimized by the outer minimization to reduce the empirical risk for adversarial samples.

The trade-off between standard and adversarial accuracy is a key factor for the real-world applications of AT (Tsipras et al., 2018; Balaji et al., 2019; Yang et al., 2020b; Stutz et al., 2019; Zhang et al., 2019). Although AT can improve robustness against adversarial samples, it also undermines the performance on standard samples. Existing AT methods (Madry et al., 2018; Cai et al., 2018; Zhang et al., 2019; Wang et al., 2019) design a hybrid loss by combining standard loss and an adversarial loss linearly, where the linear coefficient typically serves as the trade-off factor.

In this paper, we argue that linearly weighted-average method for AT, as well as the **Vanilla AT**, cannot achieve a ‘near-optimal’ trade-off. In other words, it fails to approximately achieve the Pareto optimal points on the Pareto front of standard and adversarial accuracies. We find that the conflict between the parameter gradient derived from standard loss (**standard gradient**) and the one derived from adversarial loss (**adversarial gradient**) is the main source of this failure. Such a gradient conflict causes the model parameter to be stuck in undesirable local optimal points, and it becomes more severe with the increase of adversarial attack budget. In addition, to obtain adversarial robustness, linearly weighted-average method usually sacrifices too much performance on standard samples, which hinders AT from real-world applications.

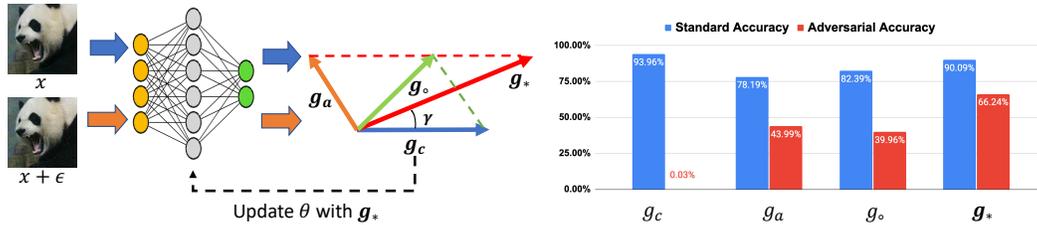


Figure 1: The key motivation of CA-AT aims to solve the conflict between clean gradient g_c and adversarial gradient g_a . Unlike the existing weighted-averaged method optimizing model parameter θ by g_o as the average of g_c and g_a (Vanilla AT), CA-AT utilizes g_* for parameter optimization by gradient projection based on a new trade-off factor ϕ . The bar chart on the right side illustrates that the model optimized by g_* (highlighted as the boldface) can achieve better standard accuracy (blue bar) and adversarial accuracy (red bar) compared to models optimized by g_o . The results of the bar chart on the right are produced by training a ResNet18 on CIFAR10 against the PGD (Madry et al., 2017) attack.

To solve the problems mentioned above, we propose **Conflict-Aware Adversarial Training (CA-AT)** to mitigate the conflict during adversarial training. Inspired by gradient surgery (Yu et al., 2020) in multi-task learning, CA-AT utilizes a new trade-off factor defined as the angle between the standard and adversarial gradients. As depicted in Fig. 1, if the angle is larger than the pre-defined trade-off factor γ , CA-AT will project the adversarial gradient onto the ‘cone’ around the standard gradient constructed based on the pre-defined trade-off factor; otherwise, it will use the standard gradient to optimize the model parameter θ directly. Compared to the linearly weighted-average AT with a fixed trade-off factor, CA-AT can boost both standard and adversarial accuracy. Our primary contributions are summarized as follows:

1. We shed light on the existence of conflict between standard and adversarial gradient which causes a sub-optimal trade-off between standard and adversarial accuracy in AT, when we optimize standard and adversarial loss in weighted-average paradigm by a fixed trade-off factor.
2. To alleviate the gradient conflict in AT, we propose a new paradigm called Conflict-Aware Adversarial Training (CA-AT). It achieves a better trade-off between standard and adversarial accuracy compared to Vanilla AT.
3. Through comprehensive experiments across a wide range of settings, we demonstrate CA-AT consistently improves the trade-off between standard and adversarial accuracy in the context of training from scratch and parameter-efficient finetuning (PEFT), across diverse adversarial loss functions, adversarial attack types, model architectures, and datasets.

2 RELATED WORKS

Adversarial Training. Adversarial training (AT) is now broadly considered as the most effective method to achieve adversarial robustness for deep learning models (Qian et al., 2022; Singh et al., 2024). The key idea of AT is to involve adversarial samples during the training process. Existing works for AT can be mainly grouped into regularization-driven and strategy-driven. For regularization-driven AT methods, the goal is to design an appropriate loss function for adversarial samples, such as cross-entropy (Madry et al., 2017), logits pairing (CLP) (Kannan et al., 2018), and TRADES (Zhang et al., 2019). On the other hand, strategy-driven AT methods focus on improving adversarial robustness by designing appropriate training strategies. For example, ensemble AT (Tramèr et al., 2017; Yang et al., 2020a) alleviates the sharp parameter curvature by utilizing adversarial examples generated from different target models, curriculum AT (Cai et al., 2018) gains adversarial robustness progressively by learning from easy adversarial samples to hard adversarial samples, and adaptive AT (Ding et al., 2018; Cheng et al., 2020; Jia et al., 2022) improves adversarial robustness by adjusting the attack intensity and attack methods. With the development of large-scale pretrained models (Kolesnikov et al., 2020; Dong et al., 2020), (Jia et al., 2024; Hua et al., 2023) demonstrates the superiority of adversarial PEFT of robust pretrained models, compared to adversarial training from scratch.

108 However, strategy-based AT methods need to involve additional attack methods or target models in
 109 the training process, which will increase the time and space complexity when we apply them. CA-AT
 110 can improve both standard and adversarial performance without any increasing cost of training time
 111 and computing resources.

112 **Gradients Operation.** Gradients Operation, also known as gradient surgery (Yu et al., 2020), aims to
 113 improve model performance by directly operating the parameter gradient during training. It was first
 114 presented in the area of multi-task learning to alleviate the gradient conflict between loss functions
 115 designed for different tasks. The conflict can be measured by cosine similarity (Yu et al., 2020) or
 116 Euclidean distance (Liu et al., 2021a) between the gradients derived from different loss functions.
 117 Besides, multi-task learning, (Mansilla et al., 2021) incorporates gradient operation to encourage
 118 gradient agreement among different source domains, enhancing the model’s generalization ability to
 119 the unseen domain, and (Chaudhry et al., 2018; Yang et al., 2023) alleviate the forgetting issue in
 120 continual learning by projecting the gradients from the current task to the orthogonal direction of
 121 gradients derived from the previous task.

122 We are the first work to observe the gradient conflict between standard and adversarial loss during AT
 123 and further reveal its relation to adversarial attack budget. Moreover, we propose a new trade-off
 124 paradigm specifically designed for AT based on gradient operation. It can achieve a better trade-off
 125 compared to Vanilla AT and guarantee the standard performance well.

127 3 GRADIENT CONFLICT IN AT

128
 129 In this section, we will discuss the occurrence of gradient conflict in AT via a synthetic dataset and
 130 real-world datasets such as CIFAR10 and CIFAR100. Additionally, we demonstrate such a conflict
 131 will become more serious with the increase of the attack budget theoretically and practically.

133 3.1 PRELIMINARIES & NOTATIONS

134
 135 Considering a set of images, each image $x \in \mathbb{R}^d$ and its label $y \in \mathbb{R}^l$ is drawn i.i.d. from distribution
 136 \mathcal{D} . The classifier $f : \mathbb{R}^d \rightarrow \mathbb{R}^l$ parameterized by θ aims to map an input image to the probabilities of
 137 the classification task. The objective of AT is to ensure that f does not only perform well on x , but
 138 also manifests robustness against adversarial perturbation ϵ bounded by attack budget δ as $\|\epsilon\|_p \leq \delta$,
 139 where p determinates the L_p norm constraint on the perturbations ϵ commonly taking on the values
 140 of ∞ or 2. The perturbation ϵ can be defined as $\epsilon = \arg \max_{\|\epsilon\|_p \leq \delta} \mathcal{L}(x + \epsilon, y; \theta)$, which can be
 141 approximated by gradient-based adversarial attacks such as PGD. Throughout the remaining part of
 142 this paper, we refer to x as the standard sample and $x + \epsilon$ as the adversarial sample.

143 We define clean loss $\mathcal{L}_c = \mathcal{L}(x, y; \theta)$ and adversarial loss $\mathcal{L}_a = \mathcal{L}(x + \epsilon, y; \theta)$, respectively. \mathcal{L} is the
 144 loss function for classification task (e.g. cross-entropy). As shown in Eq. (1), the goal of adversarial
 145 training is to obtain the parameter θ that can be both accurate and robust.

$$146 \min_{\theta} (\mathbb{E}_{(x,y) \sim \mathcal{D}}[\mathcal{L}_c], \mathbb{E}_{(x,y) \sim \mathcal{D}}[\mathcal{L}_a]) \quad (1)$$

147 For vanilla AT, as mentioned in Section 2, optimizing a hybrid loss containing standard loss \mathcal{L}_c
 148 and adversarial loss \mathcal{L}_a is a widely-used method for solving Eq. (1). As shown in Eq. (2), existing
 149 works (Wang et al., 2019; Zhang et al., 2019; Kannan et al., 2018) construct such a hybrid loss by
 150 using a linear-weighted approach for \mathcal{L}_c and \mathcal{L}_a .

$$151 \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}}[\lambda \mathcal{L}_a + (1 - \lambda) \mathcal{L}_c], \quad (2)$$

152 where $\lambda \in [0, 1]$ serves as a fixed hyper-parameter for the trade-off between \mathcal{L}_c and \mathcal{L}_a . Refer to
 153 Fig. 1, the optimization process of Eq. (2) can be described as utilizing $g_o = (1 - \lambda)g_c + \lambda g_a$ to
 154 update θ at each optimization step, where $g_c = \frac{\partial \mathcal{L}_c}{\partial \theta}$ and $g_a = \frac{\partial \mathcal{L}_a}{\partial \theta}$ represent standard and adversarial
 155 gradients, respectively.

156
 157 To measure how well we can solve Eq. (1), we define a metric $\mu = \|g_c\|_2 \cdot \|g_a\|_2 \cdot (1 - \cos(g_c, g_a))$.
 158 The basic motivation for the consideration of μ is that it should combine three kinds of signals
 159 during AT simultaneously: **(1)** $\|g_c\|_2$ reflects the convergence of clean loss \mathcal{L}_c , **(2)** $\|g_a\|_2$ reflects the
 160 convergence of adversarial loss \mathcal{L}_a , and **(3)** $(1 - \cos(g_c, g_a))$ reflects the directional conflict between
 161 g_c and g_a . Based on **(1)**, **(2)**, and **(3)**, a small μ implies that both \mathcal{L}_c and \mathcal{L}_a have converged well while
 reaching a consensus on the optimization direction for the next step.

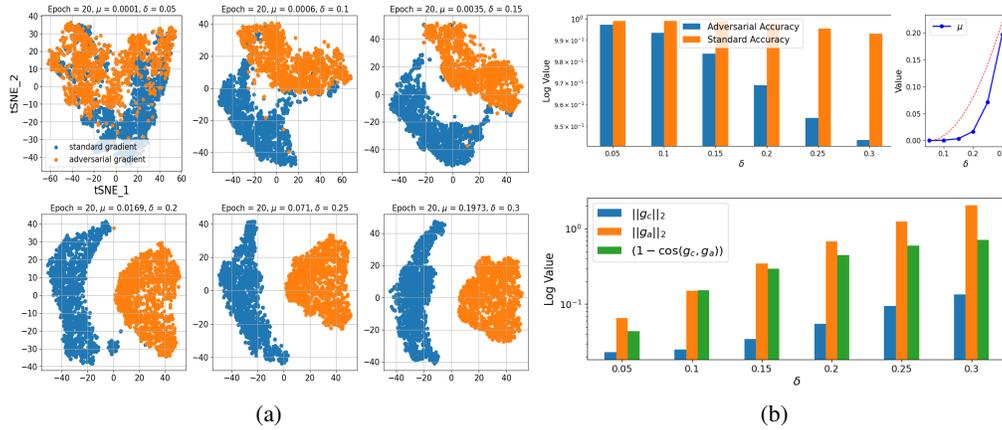


Figure 2: The experimental results of conducting Vanilla AT with $\lambda = 0.5$ for a binary classification task on our MNIST-crafted data. In Fig. 2a, each subfigure is the tSNE (Hinton & Roweis, 2002) visualization displaying the distribution of adversarial gradients (g_a) and standard gradients (g_c) for various training samples at the final epoch with different attack budgets ($\delta = [0.05, 0.1, 0.15, 0.2, 0.25, 0.3]$). In Fig. 2b, the upper bar chart shows the standard and adversarial accuracy on testing set with different δ similar to Fig. 2a. The upper left line chart shows the relation between the $\mu = \|g_c\|_2 \cdot \|g_a\|_2 \cdot (1 - \cos(g_c, g_a))$ and δ , where the red line is the theoretical upper bound presented in Theorem 1. For decomposing μ , lower bar chart shows the relation between δ and $\|g_c\|_2 \|g_a\|_2 / (1 - \cos(g_c, g_a))$, respectively.

3.2 THEORETICAL & EXPERIMENTAL SUPPORT FOR MOTIVATION

We introduce Theorem 1 that demonstrates μ can be bounded by the input dimension d and perturbation budget δ in AT.

Theorem 1. Consider the gradient conflict $\mu = \|g_c\|_2 \cdot \|g_a\|_2 \cdot (1 - \cos(g_c, g_a))$ and suppose that the input x is a d -dimensional vector:

1. Given the L_2 restriction for ϵ as $\|\epsilon\|_2 \leq \delta$, we have $\mu \leq \mathcal{O}(\delta^2)$.
2. Given the L_∞ restriction for ϵ as $\|\epsilon\|_\infty \leq \delta$, we have $\mu \leq \mathcal{O}(d^2 \delta^2)$.

The intuitive understanding of Theorem 1 is that with the increasing attack budget δ , the adversarial samples in AT will move further from the distribution \mathcal{D} of standard samples. The conflict between g_a and g_c will become more serious, and L_a and L_c will be hard to converge. Therefore, the upper bound of μ will become larger. The proof of Theorem 1 will be shown in the appendix.

Synthetic Experiment. In order to show the implications of Theorem 1 empirically, we introduce the synthetic experiment as a binary classification task by selecting digit one and digit two from MNIST with a resolution of 32×32 , and train a logistic regression model parameterized by $w \in \mathbb{R}^{(32 \times 32) \times 2}$ via BCE loss by vanilla AT for 20 epochs, where ϵ is contained by its L_∞ norm as $\|\epsilon\|_\infty \leq \delta$, and $\lambda = 0.5$ serves as the trade-off factor between standard and adversarial loss. Compared to the experiments on real-world datasets, this synthetic experiment offers a distinct advantage in terms of the ability to analytically solve the inner maximization. For real-world datasets, only numerical solutions can be derived using gradient-based attacks (e.g. PGD) during AT. These numerical solutions sometimes are not promising due to gradient masking (Athalye et al., 2018; Papernot et al., 2017). On the contrary, our synthetic experiments can ensure a high-quality solution for inner maximization, eliminating the potential effect of experimental results caused by some uncertainties such as gradient masking.

Under the circumstance of a simple logistic regression model with analytical solution for inner maximization, the hybrid loss for Vanilla AT can be presented as Eq. (3), where $\exp()$ denotes the exponential function. The details of getting the analytical solution for inner maximization will be

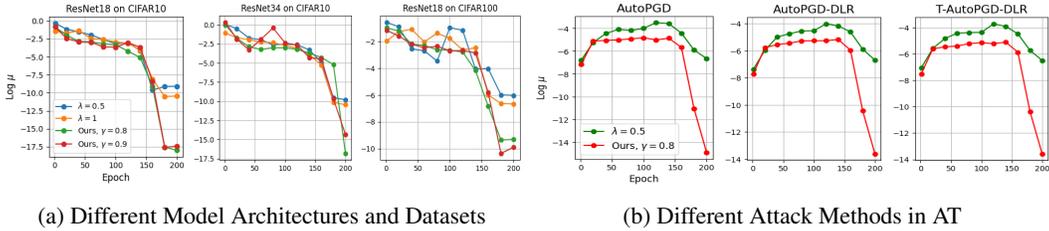


Figure 3: Results of gradient conflict metric μ on real-world datasets. Fig. 3a illustrates the results of μ among different real-world datasets (CIFAR10/CIFAR100) and model architectures (ResNet18/ResNet34), where the attack method used in AT is PGD. Fig. 3b shows the results of μ for different attack methods (AutoPGD/AutoPGD-DLR/T-AutoPGD-DLR) during AT, conducted on CIFAR10 with ResNet18.

presented in the appendix.

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\lambda \log(1 + \exp(-y \cdot (w^T x + b) + \delta \|w\|_1)) + (1 - \lambda) \log(1 + \exp(-y \cdot (w^T x + b)))] \quad (3)$$

Fig. 2 illustrates the results of this synthetic experiment. By TSNE, Fig. 2a visualizes the distributions of g_a and g_c for different training samples in the last training epoch. With the increase of attack budget δ , these two distributions are progressively fragmented, meaning g_a and g_c become more different.

For Fig. 2a, it is the tSNE visualization depicting the distributions of g_a and g_c for different training samples across varying δ . Particularly, the distributions of g_a and g_c begin to segregate more distinctly as δ becomes larger, concomitant with the increasing gradient conflict μ . Furthermore, the bar chart Fig. 2b reveals a decline in both standard and adversarial accuracies with increasing δ and μ . This trend indicates that the larger gradient conflict can harm the model’s performances on both standard and adversarial accuracies. The subfigure on the right side of Fig. 2b shows an almost quadratic growth relationship between μ and δ , the red line is the theoretical upper bound derived from Theorem 1, demonstrating the effectiveness of Theorem 1 empirically.

Experiments on Real-world Datasets. Beyond the synthetic experiment, we also conduct experiments on real-world datasets such as CIFAR10/CIFAR100, and we also observe the gradient conflict during AT. Fig. 3 shows that such a conflict exists varying from different datasets, model architectures, and attack methods used in AT, and our method ($\gamma = 0.8, \gamma = 0.9$), which will be introduced in the next section, can consistently alleviate the conflict compared to the Vanilla AT ($\lambda = 0.5, \lambda = 1$). For the fluctuation of the red line (Ours, $\gamma = 0.9$) between epochs 60–90 in the middle figure of Fig. 3a, it can be attributed to the learning rate schedule. During these epochs, the one-cycle learning rate schedule we used involves a high learning rate, which can result in increased instability and thus larger fluctuations for the gradient conflict μ .

4 METHODOLOGY

As we mentioned in Section 3, the trade-off between standard and adversarial accuracy is profoundly influenced by the gradient conflict μ (Fig. 2). The vanilla AT, which employs a linear trade-off factor λ to combine clean and adversarial loss (as seen in Eq. (2)), does not adequately address the issue of gradient conflict.

Based on this observation, we introduce Conflict-aware Adversarial Training (CA-AT) as a new trade-off paradigm for AT. The motivation of CA-AT is that the gradient conflict in AT can be alleviated by generally conducting operations on the adversarial gradient g_a and the standard gradient g_c during the training process, and such an operation should guarantee the standard accuracy because its priority is higher adversarial accuracy. Inspired by existing works related to gradient operation Yu et al. (2020); Liu et al. (2021a); Chaudhry et al. (2018); Mansilla et al. (2021), CA-AT employs a pre-defined trade-off factor γ as the goal of cosine similarity between g_c and g_a . In each iteration, instead of updating parameter θ by linearly weighted-averaged gradient g_o , CA-AT utilizes g_* to

270 update θ as Eq. (4)

$$271 \quad g_* = \begin{cases} 272 \quad g_a + \frac{\|g_a\|_2(\gamma\sqrt{1-\phi^2}-\phi\sqrt{1-\gamma^2})}{\|g_c\|_2\sqrt{1-\gamma^2}}g_c, & \phi \leq \gamma \\ 273 \quad g_c, & \phi > \gamma \end{cases} \quad (4)$$

274 where $\phi = \cos(g_a, g_c)$ is the cosine similarity between standard gradient g_c and adversarial gradient g_a . The intuitive explanation of Eq. (4) is depicted in Fig. 1. For each optimization iteration, if ϕ is less than γ , then g_* is produced by projecting g_a onto the cone of g_c at an angle $\arccos(\gamma)$. If $\phi > \gamma$, we will use the standard gradient g_c to optimize θ , because we need to guarantee standard accuracy when the conflict is not quite serious.

280 The mechanism behind Eq. (4) is straightforward. It mitigates the gradient conflict in AT by ensuring that g_c is consistently projected in a direction close to g_a . Considering an extreme case that g_c and g_a are diametrically opposite ($g_a = -g_c$), in such a scenario, if we produce the gradient by Vanilla AT as $g_o = g_c + g_a$, g_o will be a zero vector and the optimization process will be stuck. On the other hand, g_* will align closely to g_c within γ , avoiding θ to be stuck in a suboptimal point.

285 Furthermore, under the condition of $\phi \leq \gamma$, we find that CA-AT can also be viewed as a convex combination for standard and adversarial loss with a conflict-aware trade-off factor λ^* as $\mathcal{L} = \mathcal{L}_c + \lambda^*\mathcal{L}_a$, where $\lambda^* = \frac{\|g_a\|_2(\gamma\sqrt{1-\phi^2}-\phi\sqrt{1-\gamma^2})}{\|g_c\|_2\sqrt{1-\gamma^2}}$. Intuitively, λ^* increases with the decreasing of ϕ , which means we lay more emphasis on the standard loss when the conflict becomes more serious, and the hyperparameter γ here serves a role of temperature to control the intensity of changing to λ^* .

292 Algorithm 1 CA-AT

293 **Input:** Training dataset D , Loss function \mathcal{L} , Perturbation budget δ , Training epochs N , Initial model parameter θ_1 , Projection margin threshold γ , learning rate lr

294 **Output:** Trained model parameter θ_{N+1}

```

295 1: for  $t = 1$  to  $N$  do
296 2:   for each batch  $B$  in  $D$  do
297 3:      $\mathcal{L}_c = \frac{1}{|B|} \sum_{(x,y) \in B} \mathcal{L}(x, y; \theta_t)$ 
298 4:      $\mathcal{L}_a = \frac{1}{|B|} \sum_{(x,y) \in B} \max_{\|\epsilon\|_\infty \leq \delta} \mathcal{L}(x + \epsilon, y; \theta_t)$ 
299 5:      $g_c, g_a = \nabla_{\theta_t} \mathcal{L}_c, \nabla_{\theta_t} \mathcal{L}_a$ 
300 6:      $\phi = \cos(g_c, g_a)$ 
301 7:     if  $\phi < \gamma$  then
302 8:        $g_* = g_a + \frac{\|g_a\|_2(\gamma\sqrt{1-\phi^2}-\phi\sqrt{1-\gamma^2})}{\|g_c\|_2\sqrt{1-\gamma^2}}g_c$ 
303 9:     else
304 10:       $g_* = g_c$ 
305 11:    end if
306 12:     $\theta_t = \theta_t - lr * g_*$ 
307 13:  end for
308 14:   $\theta_{t+1} = \theta_t$ 
309 15: end for

```

312 The pseudo-code of the CA-AT is shown as Algorithm 1. In each training batch B , we calculate both standard loss \mathcal{L}_c and adversarial loss \mathcal{L}_a . By evaluating and adjusting the alignment between standard gradient g_c and adversarial gradient g_a , the algorithm ensures the model not only performs well via standard samples but also maintains robustness against designed perturbations. This adjustment is made by modifying the adversarial gradient g_a to better align with the standard gradient g_c based on the projection margin threshold γ , where g_* is produced to optimize the model parameter θ_t in each round t .

320 5 EXPERIMENTAL RESULTS & ANALYSIS

321 In this section, we demonstrate the effectiveness of CA-AT for achieving better trade-off results compared to Vanilla AT. We conduct experiments on adversarial training from scratch and adversarial PEFT among various datasets and model architectures. Besides, motivated by Theorem 1, we evaluate

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

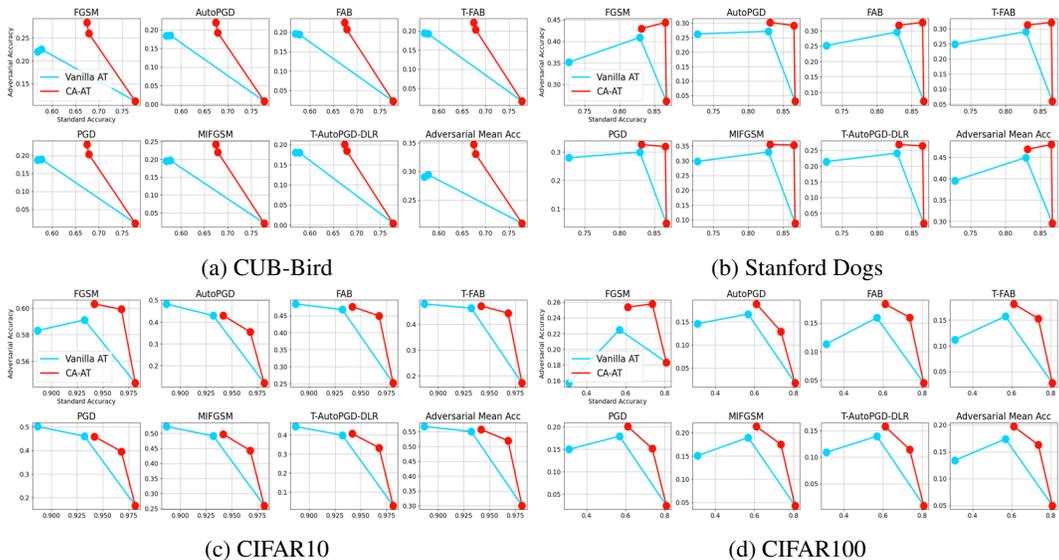


Figure 4: SA-AA Fronts for Adversarial PEFT on Swin-T using Adapter.

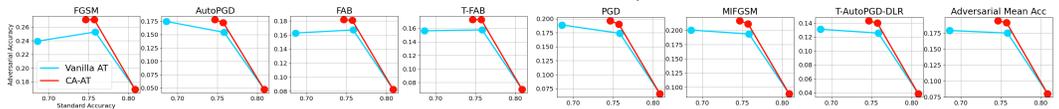


Figure 5: SA-AA Fronts for Adversarial PEFT on ViT using Adapter on Stanford Dogs.

CA-AT by involving adversarial samples with a larger budget in training. Experimental results show that CA-AT can boost the model’s robustness by handling adversarial samples with a larger budget, while Vanilla AT fails.

5.1 EXPERIMENTAL SETUP

Datasets and Models. We evaluate our proposed method on various image classification datasets including CIFAR10 (Krizhevsky et al., 2009), CIFAR100 (Krizhevsky et al., 2009), CUB-Bird (Wah et al., 2011), and StanfordDogs (Khosla et al., 2011). The model architectures we utilized to train from scratch on CIFAR10 and CIFAR100 are ResNet18, ResNet34 (He et al., 2016), and WideResNet28-10 (WRN-28-10) (Zagoruyko & Komodakis, 2016). We set the value of running mean and running variance in each Batch Normalization block into false as a trick to boost adversarial robustness (Wang et al., 2022; Walter et al., 2022). For experiments on PEFT, we fine-tune Swin Transformer (Swin-T) (Liu et al., 2021b) and Vision Transformer (ViT) (Dosovitskiy et al., 2020) by using Adapter (Pfeiffer et al., 2020; 2021), which fine-tunes the large pretrained model by inserting a small trainable module into each block. Such a module adapts the internal representations for specific tasks without altering the majority of the pretrained model’s parameters. Both Swin-T and ViT are pretrained adversarially (Dong et al., 2020) on ImageNet. For the experiments on ResNet, we set the resolution of input data as 32×32 , and use resolution as 224×224 for the PEFT experiments on Swin-T and ViT.

Hyper-parameters for AT. For adversarial training from scratch, we use the PGD attack with $\delta = 8/255$ with step size $2/255$ and step number 10. For the optimizer, we use SGD with momentum as 0.9 and the initial learning rate as 0.4. We use the one-cycle learning rate policy (Smith & Topin, 2019) as the dynamic adjustment method for the learning rate within 200 epochs. The details of hyperparameter setup for adversarial PEFT will be shown in Appendix. Generally, we use a sequence of operations as random crop, random horizontal flip, and random rotation for data augmentation. For a fair comparison, we maintain the same hyper-parameters across experiments for vanilla AT and CA-AT on both adversarial training from scratch and PEFT.

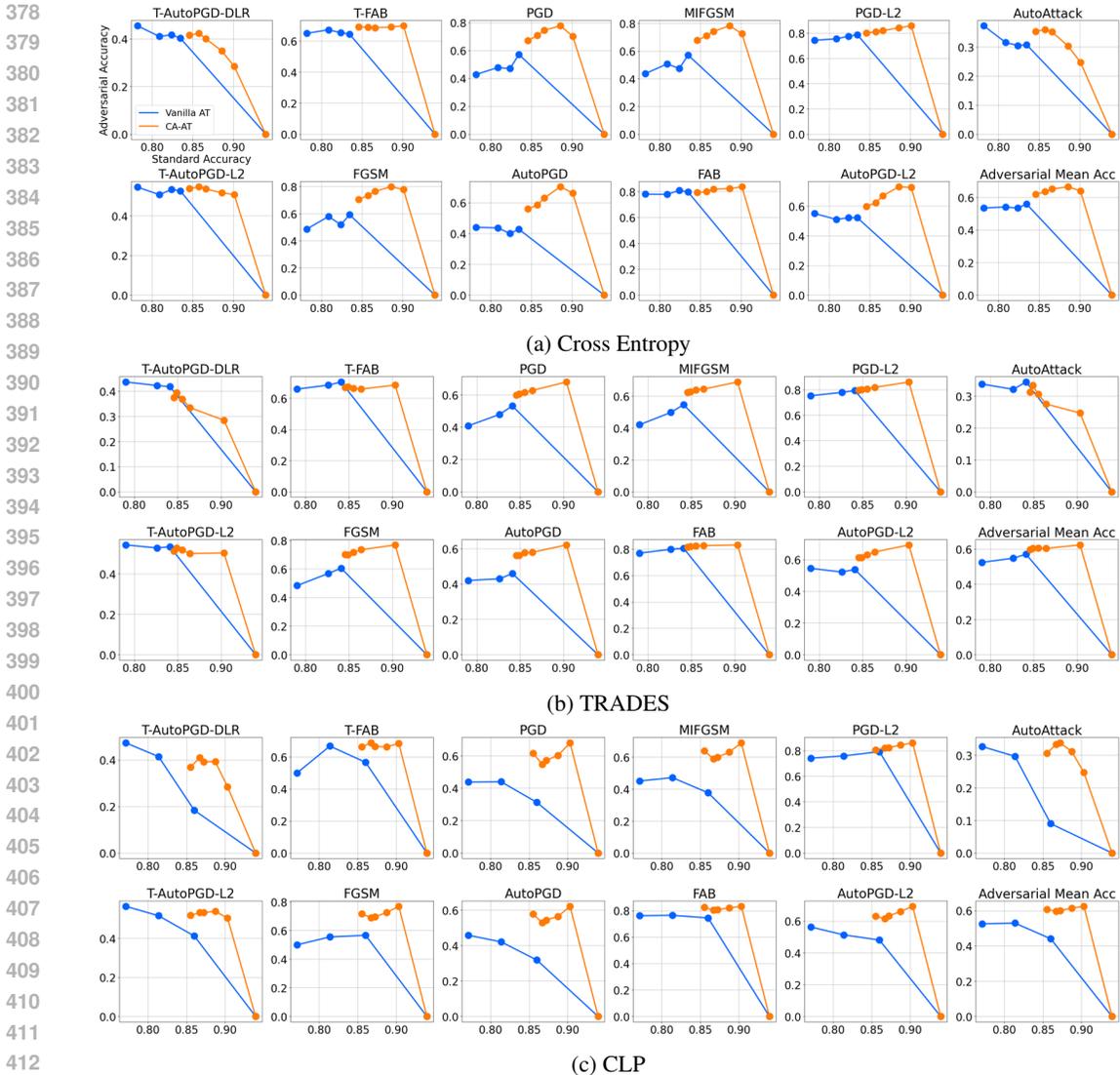


Figure 6: SA-AA Fronts for Adversarial Training from Scratch on CIFAR10 using ResNet18 with Different Adversarial Loss Functions including Cross Entropy, TRADES, and CLP.

Evaluation. We evaluate adversarial robustness by reporting the accuracies against extensive adversarial attacks constrained by L_∞ and L_2 . For attacks bounded by L_∞ norm, we selected most representative methods including PGD (Madry et al., 2018), AutoPGD (Croce, 2020), FGSM (Goodfellow et al., 2014), MIFGSM (Dong et al., 2018), FAB (Croce & Hein, 2020), and AutoAttack (Croce, 2020). Besides, we also conducted the targeted adversarial attacks, where they are denoted as a 'T-' as the prefix (e.g. T-AutoPGD). For all the targeted adversarial attacks, we set the number of classes as 10. Attacks bounded by L_2 norm are denoted as '-L2' in suffix (e.g. AutoPGD-L2). Besides, we apply attacks with different loss functions such as cross entropy (AutoPGD) and difference of logits ratio (AutoPGD-DLR), to avoid the 'fake' adversarial examples caused by gradient vanishing (Athalye et al., 2018).

To measure the quality of trade-off between standard accuracy (SA) and adversarial accuracy (AA), we define **SA-AA front** as an empirical Pareto front for SA and AA. We draw this front by conducting different λ on Vanilla AT and different γ on CA-AT.

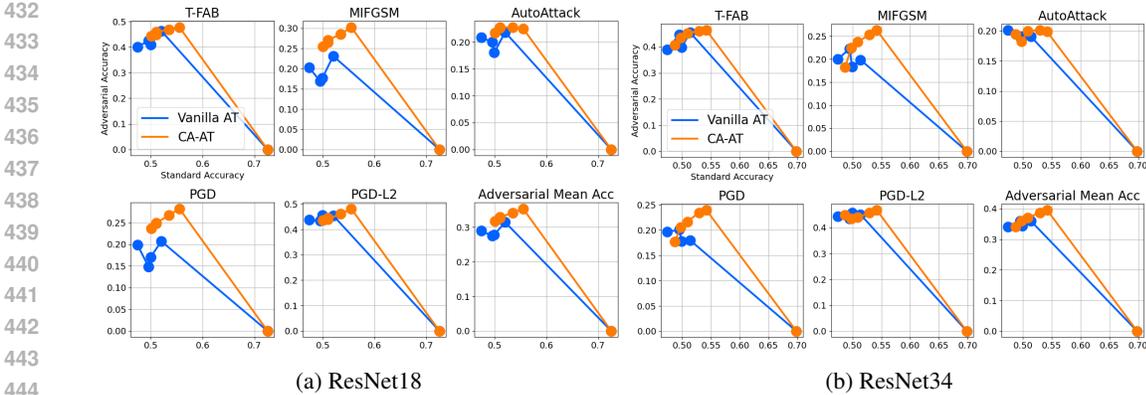


Figure 7: SA-AA Fronts for Adversarial Training from Scratch on CIFAR100.

	$p = \infty$	Standard Accuracy		PGD		AutoPGD		MIFGSM		FAB		T-FAB		FGSM	
		CA-AT	Vanilla AT	CA-AT	Vanilla AT	CA-AT	Vanilla AT	CA-AT	Vanilla AT	CA-AT	Vanilla AT	CA-AT	Vanilla AT	CA-AT	Vanilla AT
ResNet18	8/255	0.8659	0.8239	0.7442	0.4703	0.6301	0.3996	0.7419	0.4745	0.8177	0.809	0.6861	0.6538	0.7649	0.519
	16/255			0.7311	0.4248	0.5555	0.2486	0.7233	0.4225	0.7475	0.78	0.5445	0.5104	0.7435	0.4387
	24/255			0.7189	0.405	0.4886	0.1963	0.7182	0.413	0.6858	0.7333	0.4599	0.4783	0.7235	0.403
	32/255			0.7033	0.3877	0.4455	0.1589	0.7182	0.413	0.6402	0.6836	0.4044	0.4507	0.7066	0.379
ResNet34	8/255	0.8753	0.8305	0.8098	0.5973	0.7285	0.4417	0.8111	0.5983	0.8247	0.8068	0.7274	0.6951	0.8149	0.5327
	16/255			0.8034	0.5756	0.6793	0.3395	0.8077	0.5791	0.7738	0.7613	0.6013	0.5762	0.7916	0.2762
	24/255			0.7957	0.5602	0.6445	0.2859	0.8067	0.5743	0.7307	0.6937	0.5142	0.5174	0.7743	0.1428
	32/255			0.785	0.5443	0.6165	0.2498	0.8067	0.5743	0.6918	0.6221	0.4424	0.4822	0.7616	0.088

Table 1: Evaluation results on CIFAR10 for CA-AT ($\gamma = 0.8$) and Vanilla AT ($\lambda = 0.5$) across different L_∞ -based attacks with various values of budget δ .

5.2 EXPERIMENTAL RESULTS ON PEFT

CA-AT offers the better trade-off on adversarial PEFT. Fig. 4 shows the SA-AA fronts on fine-tuning robust pretrained Swin Transformer on CUB-Bird and StanfordDogs by using Adapter. We set $\lambda = [0, 0.5, 1]$ for Vanilla AT and $\gamma = [0.8, 0.9, 1]$ for CA-AT. The red data points for CA-AT are positioned in the upper right area relative to the blue points for Vanilla AT. It shows that CA-AT can consistently attain better standard and adversarial accuracy compared to the Vanilla AT across different datasets. Besides, we observed that on fine-grained datasets such as CUB-Bird and Stanford Dogs, the superiority of CA-AT is more significant compared to the results on normal datasets.

Results for CA-AT with Different Pretrained Models. Fig. 5 shows that CA-AT can also boost the trade-off performance on ViT. The main difference between these two models is that, ViT treats image patches as tokens and processes them with a standard transformer architecture Vaswani et al. (2017), while Swin-T uses shifted windows for hierarchical feature merging. While ViT applies global attention directly on image patches, Swin Transformer applies local attention within windows and uses a hierarchical approach to better handle larger and more detailed images. The superiority of CA-AT on ViT is not as significant as it is on Swin-T (Fig. 4b), but it still can gain better standard and adversarial accuracy compared to Vanilla AT.

5.3 EXPERIMENTAL RESULTS ON TRAINING FROM SCRATCH

CA-AT results in better trade-off with different adversarial loss functions. Fig. 6a visualizes SA-AA fronts from experiments using vanilla AT with $\lambda = [0, 0.25, 0.5, 0.75, 1]$ and CA-AT with $\gamma = [0.7, 0.75, 0.8, 0.85, 0.9, 1]$ on CIFAR10. In this figure, most orange data points (CA-AT) lie in the upper right space of blue points (Vanilla AT), indicating that CA-AT offers a better empirical Pareto front for the trade-off between standard accuracy and adversarial accuracy. Moreover, Fig. 6c and Fig. 6b show CA-AT can also consistently boost the adversarial accuracy for different adversarial loss functions used in AT such as TRADES (Zhang et al., 2019) and CLP (Kannan et al., 2018). For the experiments on CIFAR100, we selected the strongest and most representative attack methods to evaluate the model’s robustness, including targeted attack (T-FAB), untargeted attacks (PGD, MIFGSM), L_2 -norm attack (T-PGD), and ensemble attack (AutoAttack). Showing the trade-off

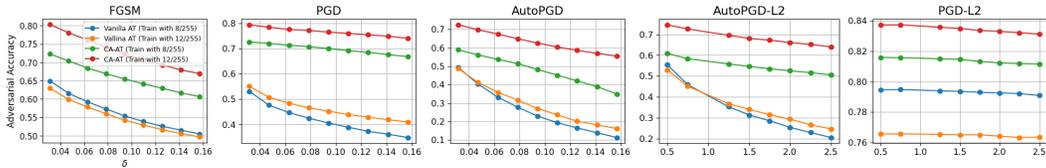


Figure 8: Results for Vanilla AT and CA-AT trained on adversarial samples with two different budget values ($\delta = 8/255, \delta = 12/255$) on CIFAR10 with ResNet18. We evaluate the adversarial accuracy among different adversarial attacks with different budget values δ denoting as the x-axis.

results on CIFAR100 in Fig. 7a (ResNet18) and Fig. 7b (ResNet34), the performance gain of CA-AT is more limited compared to the one on CIFAR10, but it can still achieve better performance on standard accuracy and adversarial accuracy against various adversarial attacks.

CA-AT is more robust to adversarial attacks with larger budget values. We evaluate adversarial precision through various adversarial attacks with different attack budget values δ , to demonstrate the superiority of our model over Vanilla AT under various intensities of adversarial attacks. We applied both Vanilla AT and CA-AT to ResNet18 on CIFAR10, and the results about L_∞ -based attacks are shown in Table 1. In Table 1, although our CA-AT achieves slightly lower adversarial accuracy against FAB when δ is larger than $8/255$, it outperforms the Vanilla AT in both standard accuracy and adversarial accuracy on any other attack methods (e.g. AutoPGD, MIFGSM, and T-FAB) with different budget δ . It clearly illustrates that, compared to Vanilla AT, CA-AT can enhance the model’s adversarial robustness ability to resist stronger adversarial attacks with larger budget δ .

CA-AT enables AT via stronger adversarial examples. In our toy experiment (Fig. 2) and Theorem 1, the conflict μ would be more serious if we utilize adversarial examples with larger attack budget δ during AT. It implies that Vanilla AT cannot handle stronger adversarial examples during training because of the gradient conflict. In Fig. 8, we visualize the results of training ResNet34 on CIFAR10 with adversarial samples produced by the same attack method (PGD), but different attack budgets ($\delta = 8/255$ and $\delta = 16/255$), and evaluate the adversarial accuracies against various adversarial methods (e.g. FGSM and PGD) with different budgets (x-axis). Compared to the blue and orange curves (Vanilla AT with $\delta = 8/255$), it shows that Vanilla AT fails when training with the adversarial attack with a higher perturbation bound, causing a decrease in both standard and adversarial accuracy. On the contrary, CA-AT, shown as the green and red curves, can improve both standard and adversarial accuracy by involving stronger adversarial samples with larger attack budgets.

Experimental Results in Appendix. More experimental results for CA-AT regarding different model architectures (WRN-28-10), different attack methods utilized for producing adversarial samples during AT, various L_2 -based attacks with different budgets, and black-box attacks can be found in Appendix C. In addition, the detailed proof for Theorem 1 is included in Appendix A.

6 CONCLUSION & OUTLOOK

In this work, we illustrate that the weighted-average method in AT is not capable of achieving the ‘near-optimal’ trade-off between standard and adversarial accuracy due to the gradient conflict existing in the training process. We demonstrate the existence of such a gradient conflict and its relation to the attack budget of adversarial samples used in AT practically and theoretically. Based on this phenomenon, we propose a new trade-off framework for AT called Conflict-Aware Adversarial Training (CA-AT) to alleviate the conflict by gradient operation. Extensive results demonstrate the effectiveness of CA-AT for gaining trade-off results under the setting of training from scratch and PEFT. **Considering the cost for gradient operation, CA-AT is more appropriate for adversarial PEFT than full fine-tuning when dealing with very large models like ViT.**

For future work, we plan to undertake a more detailed exploration of the gradient conflict phenomenon in AT from the data-centric perspective. We hold the assumption that some training samples can cause serious gradient conflict, while others do not. We will evaluate this assumption in the future work, and intend to reveal the influence of training samples causing gradient conflict.

REFERENCES

- 540
541
542 Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack:
543 a query-efficient black-box adversarial attack via random search. In *European conference on*
544 *computer vision*, pp. 484–501. Springer, 2020.
- 545 Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of
546 security: Circumventing defenses to adversarial examples. In *International conference on machine*
547 *learning*, pp. 274–283. PMLR, 2018.
- 548
549 Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. Recent advances in adversarial training
550 for adversarial robustness. *arXiv preprint arXiv:2102.01356*, 2021.
- 551 Yogesh Balaji, Tom Goldstein, and Judy Hoffman. Instance adaptive adversarial training: Improved
552 accuracy tradeoffs in neural nets. *arXiv preprint arXiv:1910.08051*, 2019.
- 553
554 Qi-Zhi Cai, Chang Liu, and Dawn Song. Curriculum adversarial training. In *Proceedings of the 27th*
555 *International Joint Conference on Artificial Intelligence*, pp. 3740–3747, 2018.
- 556
557 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017*
558 *IEEE symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.
- 559 Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopad-
560 hyay. Adversarial attacks and defences: A survey. *arXiv preprint arXiv:1810.00069*, 2018.
- 561
562 Arslan Chaudhry, Marc’ Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient
563 lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- 564
565 Minhao Cheng, Qi Lei, Pin-Yu Chen, Inderjit Dhillon, and Cho-Jui Hsieh. Cat: Customized
566 adversarial training for improved robustness. *arXiv preprint arXiv:2002.06789*, 2020.
- 567
568 Francesco Croce. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-
569 free attacks. In *International conference on machine learning*, pp. 2206–2216. PMLR, 2020.
- 570
571 Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive
572 boundary attack. In *International Conference on Machine Learning*, pp. 2196–2205. PMLR, 2020.
- 573
574 Gavin Weiguang Ding, Yash Sharma, Kry Yik Chau Lui, and Ruitong Huang. Mma training: Direct
575 input space margin maximization through adversarial training. *arXiv preprint arXiv:1812.02637*,
576 2018.
- 577
578 Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting
579 adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision*
580 *and pattern recognition*, pp. 9185–9193, 2018.
- 581
582 Yinpeng Dong, Qi-An Fu, Xiao Yang, Tianyu Pang, Hang Su, Zihao Xiao, and Jun Zhu. Benchmarking
583 adversarial robustness on image classification. In *Proceedings of the IEEE/CVF Conference on*
584 *Computer Vision and Pattern Recognition (CVPR)*, pp. 321–331, 2020.
- 585
586 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
587 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image
588 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference*
589 *on Learning Representations*, 2020.
- 590
591 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
592 examples. 2014.
- 593
594 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
595 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
596 pp. 770–778, 2016.
- 597
598 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the*
599 *IEEE international conference on computer vision*, pp. 2961–2969, 2017.

- 594 Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information*
595 *processing systems*, 15, 2002.
- 596
- 597 Andong Hua, Jindong Gu, Zhiyu Xue, Nicholas Carlini, Eric Wong, and Yao Qin. Initialization
598 matters for adversarial transfer learning. *arXiv preprint arXiv:2312.05716*, 2023.
- 599
- 600 Xiaojun Jia, Yong Zhang, Baoyuan Wu, Ke Ma, Jue Wang, and Xiaochun Cao. Las-at: adversarial
601 training with learnable attack strategy. In *Proceedings of the IEEE/CVF Conference on Computer*
602 *Vision and Pattern Recognition*, pp. 13398–13408, 2022.
- 603 Xiaojun Jia, Yong Zhang, Xingxing Wei, Baoyuan Wu, Ke Ma, Jue Wang, and Xiaochun Cao.
604 Improving fast adversarial training with prior-guided knowledge. *IEEE Transactions on Pattern*
605 *Analysis and Machine Intelligence*, 2024.
- 606
- 607 Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint*
608 *arXiv:1803.06373*, 2018.
- 609
- 610 Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for
611 fine-grained image categorization: Stanford dogs. In *Proc. CVPR workshop on fine-grained visual*
612 *categorization (FGVC)*, volume 2. Citeseer, 2011.
- 613 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and
614 Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision–ECCV*
615 *2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pp.
616 491–507. Springer, 2020.
- 617 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 618
- 619 Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for
620 multi-task learning. *Advances in Neural Information Processing Systems*, 34:18878–18890, 2021a.
- 621 Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples
622 and black-box attacks. In *International Conference on Learning Representations*, 2016.
- 623
- 624 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
625 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
626 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021b.
- 627
- 628 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
629 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
630 2017.
- 631 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
632 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
633 *Learning Representations*, 2018.
- 634
- 635 Lucas Mansilla, Rodrigo Echeveste, Diego H Milone, and Enzo Ferrante. Domain generalization via
636 gradient surgery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
637 pp. 6630–6638, 2021.
- 638
- 639 Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and
640 accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer*
641 *vision and pattern recognition*, pp. 2574–2582, 2016.
- 642
- 643 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
644 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
645 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
646 27730–27744, 2022.
- 647
- 648 Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram
649 Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on*
650 *Asia conference on computer and communications security*, pp. 506–519, 2017.

- 648 Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulic, Sebastian Ruder,
649 Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers.
650 *EMNLP 2020*, pp. 46, 2020.
- 651
652 Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfu-
653 sion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference*
654 *of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp.
655 487–503, 2021.
- 656 Zhuang Qian, Kaizhu Huang, Qiu-Feng Wang, and Xu-Yao Zhang. A survey of robust adversarial
657 training in pattern recognition: Fundamental, theory, and methodologies. *Pattern Recognition*, 131:
658 108889, 2022.
- 659 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
660 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
661 models from natural language supervision. In *International conference on machine learning*, pp.
662 8748–8763. PMLR, 2021.
- 663
664 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
665 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
666 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 667 Jérôme Rony, Luiz G Hafemann, Luiz S Oliveira, Ismail Ben Ayed, Robert Sabourin, and Eric
668 Granger. Decoupling direction and norm for efficient gradient-based l2 adversarial attacks and
669 defenses. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
670 pp. 4322–4330, 2019.
- 671
672 Naman Deep Singh, Francesco Croce, and Matthias Hein. Revisiting adversarial training for imagenet:
673 Architectures, training and generalization across threat models. *Advances in Neural Information*
674 *Processing Systems*, 36, 2024.
- 675 Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using
676 large learning rates. In *Artificial intelligence and machine learning for multi-domain operations*
677 *applications*, volume 11006, pp. 369–386. SPIE, 2019.
- 678
679 David Stutz, Matthias Hein, and Bernt Schiele. Disentangling adversarial robustness and generaliza-
680 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
681 pp. 6976–6987, 2019.
- 682 Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick Mc-
683 Daniel. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*,
684 2017.
- 685
686 Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry.
687 Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- 688 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
689 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
690 *systems*, 30, 2017.
- 691
692 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd
693 birds-200-2011 dataset. 2011.
- 694
695 Nils Philipp Walter, David Stutz, and Bernt Schiele. On fragile features and batch normalization in
696 adversarial training. *arXiv preprint arXiv:2204.12393*, 2022.
- 697
698 Haotao Wang, Aston Zhang, Shuai Zheng, Xingjian Shi, Mu Li, and Zhangyang Wang. Removing
699 batch normalization boosts adversarial training. In *International Conference on Machine Learning*,
700 pp. 23433–23445. PMLR, 2022.
- 701
702 Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving
adversarial robustness requires revisiting misclassified examples. In *International conference on
learning representations*, 2019.

702 Huanrui Yang, Jingyang Zhang, Hongliang Dong, Nathan Inkawhich, Andrew Gardner, Andrew
703 Touchet, Wesley Wilkes, Heath Berry, and Hai Li. Dverge: diversifying vulnerabilities for
704 enhanced robust generation of ensembles. *Advances in Neural Information Processing Systems*,
705 33:5505–5515, 2020a.

706 Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Russ R Salakhutdinov, and Kamalika Chaud-
707 huri. A closer look at accuracy vs. robustness. *Advances in neural information processing systems*,
708 33:8588–8601, 2020b.

709 Zeyuan Yang, Zonghan Yang, Peng Li, and Yang Liu. Restricted orthogonal gradient projection for
710 continual learning. *arXiv preprint arXiv:2301.12131*, 2023.

711 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
712 Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*, 33:
713 5824–5836, 2020.

714 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision*
715 *Conference 2016*. British Machine Vision Association, 2016.

716 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
717 Theoretically principled trade-off between robustness and accuracy. In *International conference on*
718 *machine learning*, pp. 7472–7482. PMLR, 2019.

719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755