

# Fast Adversarial Training with Adaptive Step Size

Zhichao Huang<sup>1</sup>, Yanbo Fan<sup>2</sup>, Chen Liu<sup>3</sup>, Weizhong Zhang<sup>1</sup>, Yong Zhang<sup>2</sup>

Mathieu Salzmann<sup>3</sup>, Sabine Süsstrunk<sup>3</sup>, Jue Wang<sup>2</sup>

<sup>1</sup>The Hong Kong University of Science and Technology      <sup>2</sup>Tencent AI Lab

<sup>3</sup>École Polytechnique Fédérale de Lausanne

## Abstract

While adversarial training and its variants have shown to be the most effective algorithms to defend against adversarial attacks, their extremely slow training process makes it hard to scale to large datasets like ImageNet. The key idea of recent works to accelerate adversarial training is to substitute multi-step attacks (e.g., PGD) with single-step attacks (e.g., FGSM). However, these single-step methods suffer from catastrophic overfitting, where the accuracy against PGD attack suddenly drops to nearly 0% during training, destroying the robustness of the networks. In this work, we study the phenomenon from the perspective of training instances. We show that catastrophic overfitting is instance-dependent and fitting instances with larger gradient norm is more likely to cause catastrophic overfitting. Based on our findings, we propose a simple but effective method, *Adversarial Training with Adaptive Step size (ATAS)*. ATAS learns an instance-wise adaptive step size that is inversely proportional to its gradient norm. The theoretical analysis shows that ATAS converges faster than the commonly adopted non-adaptive counterparts. Empirically, ATAS consistently mitigates catastrophic overfitting and achieves higher robust accuracy on CIFAR10, CIFAR100 and ImageNet when evaluated on various adversarial budgets.

## 1 Introduction

Adversarial examples [27] cause serious safety concerns in deploying deep learning models. In order to defend against adversarial attacks, many approaches have been proposed [10, 16, 19, 33]. Among them, adversarial training and its variants [19, 29, 33] have been recognized as the most effective defense mechanism. Adversarial training (AT) is generally formulated as a minimax problem

$$\min_{\boldsymbol{\theta}} \max_{\mathbf{x}_i^* \in \mathcal{B}_p(\mathbf{x}_i, \varepsilon)} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^*, y_i; \boldsymbol{\theta}), \quad (1)$$

where  $\mathcal{D} = (\mathbf{x}_i, y_i)_{i=1}^n$  is the training set and  $\ell(\mathbf{x}, y; \boldsymbol{\theta})$  is the loss function parametrized by  $\boldsymbol{\theta}$ .  $\mathcal{B}_p(\mathbf{x}_i, \varepsilon)$  represents a  $L_p$  norm ball centered at  $\mathbf{x}_i$  with radius  $\varepsilon$ . AT in Equation (1) boosts the adversarial robustness by adopting adversarial examples generated in the inner maximization.

Despite the effectiveness of AT, solving the inner maximization requires multiple steps of projected gradient descent (PGD) [19, 23]. Therefore, AT is much slower than vanilla training (e.g., 10 times longer training time for AT in [23]), making it challenging to scale AT to large datasets such as ImageNet.

Currently, the typical solution to accelerate AT is to substitute multi-step attacks (e.g., PGD) with single-step attacks (e.g., FGSM). Several works have been proposed following this direction, including FGSM-RS [30], ATTA [34] etc. These methods achieve the best robust accuracy for fast AT. However, recent works [2, 13] demonstrate that the single-step method suffers from catastrophic overfitting, where the model’s robustness against PGD attack suddenly drops to nearly 0% while the robust accuracy against FGSM attack rapidly increases [30]. This will completely destroy the robustness of the networks. It is worth noting that catastrophic overfitting is different from robust overfitting mentioned in [23]. The latter one refers to the generalization gap between training and test data while catastrophic overfitting means the overfitting to a specific type of attack that is irrelevant to the training and test set. Some works have been proposed to understand and alleviate the catastrophic overfitting [2, 13]. However, their solutions significantly increase the training time. For example, the gradient align regularizer  $\mathbb{E}_{\boldsymbol{\eta} \sim \mathcal{U}([- \varepsilon, \varepsilon]^d)} [1 - \cos(\nabla_{\mathbf{x}} \ell(\mathbf{x}, y; \boldsymbol{\theta}), \nabla_{\mathbf{x}} \ell(\mathbf{x} + \boldsymbol{\eta}, y; \boldsymbol{\theta}))]$  in [2] requires calculating the second order gradient and it is still 5 times slower than vanilla training. And [13] needs to check several points within the  $\ell_p$  norm ball, which needs several forward propagation and is still about 4 times slower than vanilla

training. Therefore, existing methods are still unsatisfactory in terms of both training efficiency and robust performance.

In this work, we analyze catastrophic overfitting from the perspective of training instances. By taking the gradient norm as an indicator, we find that different training instances have different probabilities of causing catastrophic overfitting. Instances with large gradient norm are more sensitive to the adversarial noise and their loss landscape is less smooth. Thus, fitting them with FGSM is more likely to distort the loss landscape, resulting in catastrophic overfitting.

Furthermore, catastrophic overfitting is closely related to the optimization process of the inner maximization, *e.g.*, the setting of step size. When catastrophic overfitting does not occur, the larger step size leads to a stronger attack and thus strengthens the robustness of the network [30]. On the other side, a larger step size is more likely to cause catastrophic overfitting in the training process [2, 30]. Based on these findings, we propose *Adversarial Training with Adaptive Step size (ATAS)*, an simple but effective fast AT method that uses the previous initialization in ATTA [34] and takes the step size of the inner maximization inversely proportional to the input gradient norm. Instances with large gradient norm are given a small step size to prevent catastrophic overfitting. By contrast, instances with small gradient norms will have large step sizes to improve the strength of the attack.

We theoretically analyze the convergence of ATAS and prove that it converges faster than the non-adaptive counterpart, which is commonly adopted in existing works [34], especially when the distribution of the input gradient norm is long-tailed. Empirically, We evaluate ATAS on CIFAR10, CIFAR100 [15] and ImageNet [6] with different network architectures and adversarial budgets, showing that ATAS mitigates catastrophic overfitting and achieves higher robust accuracy under various attacks including PGD10, PGD50 [19] and AutoAttack [5].

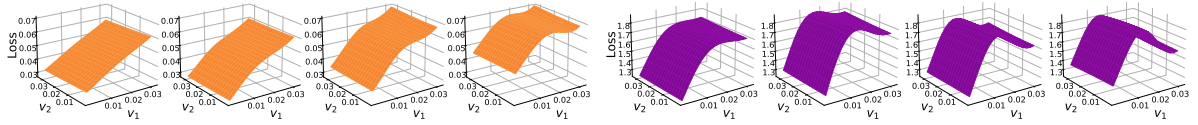
Our contributions are summarized as follows: 1) To the best of our knowledge, we are the first to analyze catastrophic overfitting from the perspective of training instances, and demonstrate that instances with large input gradient norms are more likely to cause catastrophic overfitting. 2) Based on our findings, we propose a new algorithm, ATAS, which takes the step size of the inner maximization to be inversely proportional to the input gradient norm in order to prevent catastrophic overfitting and maintain the strength of the attack. 3) Theoretically, we prove that ATAS converges faster than its non-adaptive counterpart. 4) Empirically, we conduct extensive experiments to evaluate ATAS on different datasets, network architectures and adversarial budgets, showing that ATAS consistently improves the robust accuracy and mitigates catastrophic overfitting.

## 2 Background and Related Work

### 2.1 Adversarial Examples.

Adversarial examples are first discussed in [27], where a small perturbation of the input significantly changes the prediction. Adversarial examples can be generated using the gradient of the input  $\mathbf{x}$ . Fast Gradient Signed Method (FGSM) [9] approximates the loss function  $\ell(\mathbf{x}, y; \boldsymbol{\theta})$  with the first order Taylor expansion so that adversarial examples can be generated with one step of projected gradient  $\mathbf{x}^{\text{FGSM}} = \mathbf{x} + \varepsilon \cdot \text{sgn}(\nabla_{\mathbf{x}} \ell(\mathbf{x}, y; \boldsymbol{\theta}))$ , where  $\varepsilon$  is the adversarial budget. Projected Gradient Descent (PGD) [19] extends FGSM to multiple steps to strengthen the attack. With a step size  $\alpha$ , the adversarial example at the  $t$ -th step is  $\mathbf{x}^{t+1} = \Pi_{\mathcal{B}_p(\mathbf{x}, \varepsilon)}[\mathbf{x}^t + \alpha \cdot \text{sgn}(\nabla_{\mathbf{x}^t} \ell(\mathbf{x}^t, y; \boldsymbol{\theta}))]$ , where  $\Pi_{\mathcal{B}_p(\mathbf{x}, \varepsilon)}$  means the projection onto  $\mathcal{B}_p(\mathbf{x}, \varepsilon)$ . Several stronger attacks are proposed to reliably evaluate the models' robustness [1, 4, 5]. Among them, Autoattack [5] stands out as the strongest attack.

While many algorithms [10, 16, 19, 26, 29, 33] have been proposed to defend against adversarial attacks, adversarial training and its variants [19, 29, 33] are shown to be the most effective methods to train a truly robust network. Adversarial training can be formulated as a minimax problem in Equation (1). Finding solutions of the minimax optimization has been a major endeavor in mathematics and computer science [3, 24]. Theoretically, the well-known Stochastic Gradient Descent Ascent (SGDA) algorithm finds an  $\varepsilon$ -approximate stationary point in  $\mathcal{O}(1/\varepsilon^2)$  iterations with averaging for convex-concave games [20]. However, it is not appropriate to formulate the optimization of AT as SGDA or SGDmax [17], since it only updates a part of the coordinates in  $\mathbf{x} = [x_1, x_2, \dots, x_n]$  for the maximization. The inner maximization actually corresponds to the stochastic block coordinate ascent. Empirically, the neural network is non-concave with respect to the input, so perfectly solving the inner maximization is NP-hard. It is usually approximated by a strong attack like PGD [19], which needs multiple steps of the calculation the gradients. Therefore, adversarial training is much slower than vanilla training.



(a)  $\mathcal{D}_1^1$  (Instances with smallest 10% gradient norm)    (b)  $\mathcal{D}_{10}^{10}$  (Instances with largest 10% gradient norm)

Figure 1: The loss surface of the subsets  $\mathcal{D}_1^1$  and  $\mathcal{D}_{10}^{10}$ . We average the loss of the instances from each subset.  $v_1$  is the direction of adversarial noise and  $v_2$  is a random direction. Figures from left to right plot the loss surface as the training step increases and each column of (a) and (b) corresponds to the same step of FGSM-RS.

## 2.2 Fast Adversarial Training.

FreeAT [25] first proposes a fast AT method by simultaneously optimizing the model’s parameter and the adversarial perturbations by batch replaying. YOPO [32] adopts a similar strategy to optimize the adversarial loss function. Later on, single-step methods are shown to be more effective than FreeAT and YOPO [30]. FGSM with Random Start (FGSM-RS) can be used to generate adversarial perturbations in one step to train a robust network if the hyperparameters are carefully tuned [30]. ATTA [34] utilizes the transferability of adversarial examples between epochs, using adversarial example of the previous epoch as the initialization, optimizing the model parameters with

$$\begin{aligned} \mathbf{x}_i^j &= \Pi_{\mathcal{B}_p(\mathbf{x}_i, \varepsilon)}[\mathbf{x}_i^{j-1} + \alpha \cdot \text{sgn}(\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y; \boldsymbol{\theta}))] \\ \boldsymbol{\theta} &= \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}_i^j, y; \boldsymbol{\theta}), \end{aligned} \quad (2)$$

where  $\mathbf{x}_i^j$  means the adversarial examples generated for the  $i$ -th instance  $\mathbf{x}_i$  at the  $j$ -th epoch. ATTA shows comparable robust accuracy with FGSM-RS. SLAT [21] perturbs both inputs and the latents simultaneously with FGSM, ensuring more reliable performance.

As mentioned above, these single-step methods suffer from *catastrophic overfitting*, meaning the robustness against PGD attack suddenly drops to nearly 0% while the robust accuracy against FGSM attack rapidly increases. In order to prevent catastrophic overfitting, FGSM-GA [2] adds a regularizer that aligns the direction of the input gradient. Another work [13] studies the phenomenon from the perspective of loss landscape, finding that catastrophic overfitting is a result of highly distorted loss surface. It proposes a new algorithm to resolve catastrophic overfitting by checking the loss value along the direction of the gradient. However, both algorithms require much more computation than FGSM-RS [30] and ATTA [34]. Compared with these works, we study catastrophic overfitting from the perspective of training instances and show that using adaptive step sizes in single-step methods prevents catastrophic overfitting. Our method achieves better performance with negligible computational overhead. Adaptive step sizes have been widely used in training neural networks such as AdaGrad [7], RMSProp [28] and ADAM [8, 14, 22]. However, our motivation is different, and to the best of our knowledge, we are the first to introduce the adaptive step size in fast AT.

## 3 Motivation

Catastrophic overfitting is interpreted as a result of highly distorted loss landscapes of the input [13]. For examples, FGSM-RS [30] uses large step sizes in the inner maximization to generate adversarial examples. It may only minimize the classification loss near the boundary of the adversarial budget, while the loss inside the adversarial budget may increase, leading to a highly distorted loss landscapes.

Recalling that different inputs have different loss landscape, they may result in different probabilities of causing catastrophic overfitting. Instances with large gradient norms are more sensitive to the adversarial noise. Thus, the network may simply minimize the loss on the FGSM-perturbed examples near the boundary instead of the whole space within the adversarial budget. This leads to highly distorted loss landscapes and catastrophic overfitting. The following experiments verify our hypothesis of catastrophic overfitting in FGSM-RS. The results of ATTA [34] are deferred to the Appendix B.1.

**Metrics of Input Gradient Norm.** To verify the hypothesis that instances with large gradient norms cause catastrophic overfitting, we divide the training instances into different subsets according to their gradient norms. Following the grouping method in [18], we also average the gradient norm across the training process to reduce the randomness. Formally speaking, we perform FGSM-RS to train a ResNet-18 (RN-18) on CIFAR10 for  $N = 30$  epochs with  $\varepsilon = 8/255$  and step size  $\alpha = 10/255$ . And catastrophic

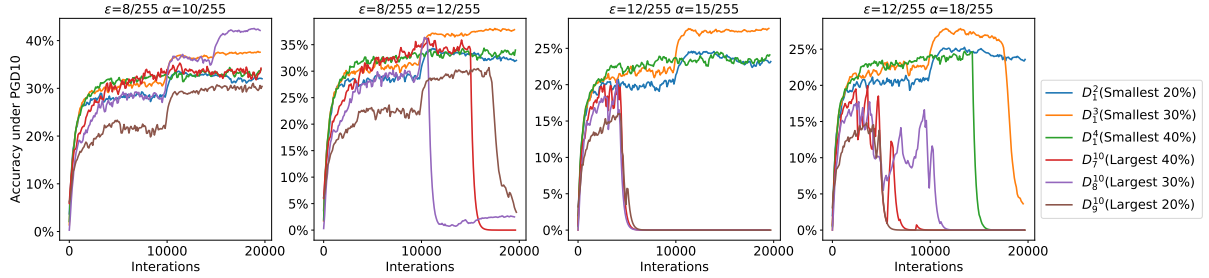


Figure 2: The robust training accuracy curve of FGSM-RS trained on different subsets of CIFAR10. The adversarial budgets and the step sizes are shown on top of each figure. The sudden decrease in accuracy indicates catastrophic overfitting.

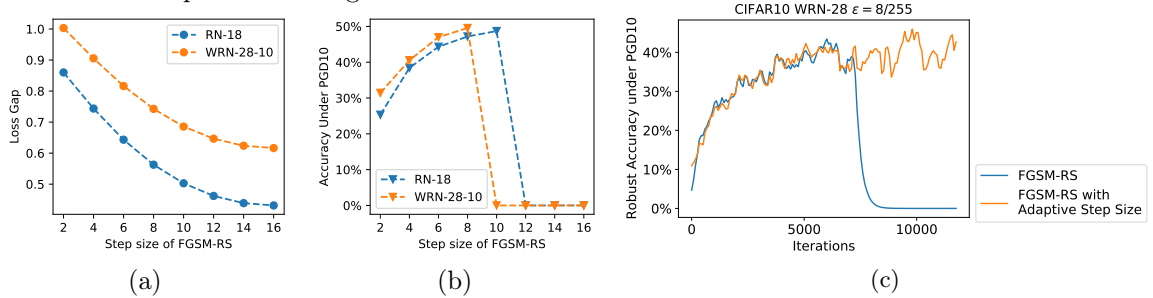


Figure 3: (a) The loss gap of training instances between PGD10 and FGSM-RS  $\ell(\mathbf{x}^{\text{PGD}}, y) - \ell(\mathbf{x}^{\text{FGSM-RS}}, y)$  with different step sizes for a FGSM-RS trained robust model; (b) The test robust accuracy of the models trained by FGSM-RS with different step sizes. (c) Accuracy of a WideResNet-28-10 under PGD10 of FGSM-RS and FGSM-RS with adaptive step size.

overfitting does not happen in this case. The average gradient norm  $GN(\mathbf{x}_i) = \frac{1}{N} \sum_{j=1}^N \|\nabla_{\tilde{\mathbf{x}}_i^j} \ell(\tilde{\mathbf{x}}_i^j, y_i; \theta)\|_2$ , where  $\tilde{\mathbf{x}}_i^j$  is the random initialization of  $\mathbf{x}_i$  at the  $j$ -th epoch. We sort  $\mathbf{x}_i$  according to  $GN(\mathbf{x}_i)$  and define  $\text{rank}(\mathbf{x}_i) = \frac{1}{n} \sum_{j=1}^n 1(GN(\mathbf{x}_j) < GN(\mathbf{x}_i))$  as the fraction of instances with smaller average gradient norm than  $\mathbf{x}_i$ . We divide the subsets according to  $\text{rank}(\mathbf{x}_i)$ :  $\mathcal{D}_i^j = \{\mathbf{x}_k | \frac{10(i-1)}{n} \leq \text{rank}(\mathbf{x}_k) < \frac{10j}{n}\}$ . The classes of each subset is balanced. The maximum and minimum proportion of one class in all subsets is 10.86% and 8.98% in CIFAR10.

**Loss Landscape.** We train a new RN-18 using FGSM-RS and enlarge the step size to  $\alpha = 14/255$  to cause catastrophic overfitting. Figure 1 shows the loss surface of the subsets with the smallest ( $\mathcal{D}_1^1$ ) and the largest gradient norm ( $\mathcal{D}_{10}^{10}$ ) when the catastrophic overfitting happens.  $\mathcal{D}_{10}^{10}$  first exhibits the catastrophic overfitting, where the loss surface of the input gets highly distorted and the loss function reaches its highest value in the middle of the adversarial budget. By contrast, the loss surface of  $\mathcal{D}_1^1$  is less distorted. Figure 1 infers that the subsets with large gradient norm are more likely to suffer from catastrophic overfitting.

**Training with Different Subsets.** We perform FGSM-RS on different subsets of CIFAR10 with different adversarial budgets  $\epsilon$  and step size  $\alpha$  to show that fitting examples with larger gradient norm is more likely to cause catastrophic overfitting. We train the RN-18 on instances with small gradient norm  $\mathcal{D}_1^2, \mathcal{D}_1^3, \mathcal{D}_1^4$  and instances with large gradient norm  $\mathcal{D}_7^{10}, \mathcal{D}_8^{10}, \mathcal{D}_9^{10}$ . While different subsets contain different number of instances, we keep the number of the training iterations the same for fair comparison. In Figure 2, we show the robust accuracy of the whole training set under PGD-10. For  $\epsilon = 8/255$  with  $\alpha = 10/255$ , the models trained with all subsets do not exhibit catastrophic overfitting. However, as the step size  $\alpha$  increases, subsets with large norms first exhibit catastrophic overfitting, while catastrophic overfitting is less likely to occur in the model trained with the subsets of small gradient norm. The figure shows 1) for each subset, catastrophic overfitting is more likely to occur when increasing the step size; 2) for a fixed step size, catastrophic overfitting is less likely to happen for subset with small gradient norm.

## 4 Algorithms

From our analysis in Section 3, the step size of the inner maximization plays an important role for the performance of the single step methods. Overly large step size draws all FGSM-perturbed noise near the boundary, causing catastrophic overfitting and thus the robust accuracy under PGD decreases to zero. However, we cannot simply reduce the step size. As shown in Figure 3a and 3b, increasing step size can

---

**Algorithm 1** ATAS

---

**Input:** Training set  $\mathcal{D}$ , The model  $f_{\theta}$  with loss function  $\ell$ , Adversarial budget  $\varepsilon$

**Output:** Optimized model  $f_{\theta^*}$

- 1:  $v_i^0 = 0$  for  $i = 1, \dots, n$
  - 2:  $\mathbf{x}_i^0 = \mathbf{x}_i + \text{Uniform}(-\varepsilon, \varepsilon)$  for  $i = 1, \dots, n$
  - 3: **for**  $j = 1$  to  $N$  **do**
  - 4:   **for**  $\mathbf{x}_i, y_i \in \mathcal{D}$  **do**
  - 5:      $v_i^j = \beta v_i^{j-1} + (1 - \beta) \|\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y_i; \theta)\|_2^2$
  - 6:      $\alpha_i^j = \gamma / (c + \sqrt{v_i^j})$
  - 7:      $\mathbf{x}_i^j = \Pi_{\mathcal{B}_p(\mathbf{x}_i, \varepsilon)}[\mathbf{x}_i^{j-1} + \alpha_i^j \cdot \text{sgn}(\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y; \theta))]$
  - 8:      $\theta = \theta - \eta \nabla_{\theta} \ell(\mathbf{x}_i^j, y; \theta)$
  - 9:   **end for**
  - 10: **end for**
- 

strengthens the adversarial attack and improves the robust accuracy.

To strengthen the attack as much as possible as well as avoid the catastrophic overfitting, we advocate utilizing the instance-wise step-size. The analysis in Section 3 shows that we should use small step sizes for instances with large gradient norms to prevent catastrophic overfitting, and large step sizes for instances with small gradient norms to the strengthen the attack. Thus, we use the moving average of the gradient norm

$$v_i^j = \beta v_i^{j-1} + (1 - \beta) \|\nabla_{\tilde{\mathbf{x}}_i} \ell(\tilde{\mathbf{x}}_i, y_i; \theta)\|_2^2, \quad (3)$$

to adjust the step size  $\alpha_i^j$  for the  $\mathbf{x}_i$  at the  $j$ -th epoch. Here,  $\tilde{\mathbf{x}}_i$  is the initialization of  $\mathbf{x}_i$  and  $\beta$  is the momentum factor stabilizing the step size. The step size  $\alpha_i^j$  is inversely proportional to  $v_i^j$ :

$$\alpha_i^j = \gamma / (c + \sqrt{v_i^j}), \quad (4)$$

where  $\gamma$  is a pre-defined learning rate and  $c$  is a constant preventing  $\alpha_i^j$  from being too large. We incorporate the adaptive step size  $\alpha_i^j$  with FGSM-RS, which randomly initializes the perturbation at the inner maximization step. The results are shown in Figure 3c, where the catastrophic overfitting does not occur by adaptive step size. In addition, the average step size of the adaptive step size method is  $10.8/255$ , which is even larger than the fixed step size  $8/255$  in FGSM-RS, leading to a stronger attack and better adversarial robustness.

Random initialization limits the magnitude of perturbations for instances with small step size, weakening the attack strength. In order to make the whole space within the adversarial budget reachable, we consider the previous initialization in ATTA [34], which utilizes the transferability of adversarial examples and uses the adversarial perturbation obtained in the previous epoch as the initialization for the inner maximization. Combined with the previous initialization, ATAS does not need large  $\alpha_i^j$  to reach the whole  $\ell_p$  norm ball. For each instance, we use adaptive step size  $\alpha_i^j$  and perform the following inner maximization to obtain the adversarial examples:

$$\mathbf{x}_i^j = \Pi_{\mathcal{B}_p(\mathbf{x}_i, \varepsilon)}[\mathbf{x}_i^{j-1} + \alpha_i^j \cdot \text{sgn}(\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y_i; \theta))], \quad (5)$$

where  $\mathbf{x}_i^j$  is the adversarial example at the  $j$ -th epoch. Then the parameter  $\theta$  is updated with  $\mathbf{x}_i^j$

$$\theta = \theta - \eta \nabla_{\theta} \ell(\mathbf{x}_i^j, y_i; \theta). \quad (6)$$

In contrast to previous methods [2, 13] that needs large computational overhead to resolve the problem of catastrophic overfitting, the overhead of ATAS is negligible, since the input gradient  $\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y_i; \theta)$  is already calculated in the attack step in Equation (5). Thus, calculating the pre-conditioner  $v_i^j$  and the step size  $\alpha_i^j$  does not need additional forward-backward passes of the network. The training time of ATAS is almost the same as ATTA [34] and FGSM-RS [30]. The detailed algorithm of ATAS is shown in Algorithm 1.

**Theoretical Analysis of ATAS.** We analyze the convergence of ATAS with  $L_{\infty}$  adversarial budget. The proof is deferred to Appendix A. Given the objective function

$$\phi(\theta, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, y_i; \theta), \quad (7)$$

the minimax problem can be formulated as follows:

$$\min_{\boldsymbol{\theta}} \max_{\mathbf{x}^* = [\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_n^*] \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}, \mathbf{x}^*), \quad (8)$$

where  $\mathbf{x}^*$  is the optimal adversarial example depending on  $\boldsymbol{\theta}$ . We consider the minimax optimization in convex-concave and smooth setting. And the loss function  $\ell$  satisfies the following assumptions.

**Assumption 4.1.** *The training loss function  $\ell$  satisfies the following constraints:*

1.  $\ell$  in convex and  $L_{\theta}$ -smooth in  $\boldsymbol{\theta}$ ;  $\boldsymbol{\theta}$  and the gradient of  $\boldsymbol{\theta}$  are bounded in the  $L_2$  norm balls

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq D_{\theta,2}, \quad \frac{1}{n} \sum_{i=1}^n \|\nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}'_i, y_i; \boldsymbol{\theta})\|_2^2 \leq G_{\theta,2}^2,$$

where  $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}, \mathbf{x}^*)$ .

2.  $\ell$  in concave and  $L_x$ -smooth in each  $\mathbf{x}_i$ .  $\mathbf{x}_i \in \mathbb{R}^d$  is bounded in an  $L_{\infty}$  norm ball with  $D_{x,\infty} = 2\varepsilon$ . For any  $\mathbf{x}$  and  $\mathbf{x}'$ ,  $\|\mathbf{x} - \mathbf{x}'\|_{\infty} \leq D_{x,\infty}$ , and the gradients of the inputs also satisfy

$$\|\nabla_{\mathbf{x}'_i} \ell(\mathbf{x}'_i, y_i; \boldsymbol{\theta})\|_2^2 \leq G_{x_i,2}^2, \quad \sum_{i=1}^n \|\nabla_{\mathbf{x}'_i} \ell(\mathbf{x}'_i, y_i; \boldsymbol{\theta})\|_2^2 \leq G_{x,2}^2$$

We average the trajectory of  $T$ -steps  $\bar{\boldsymbol{\theta}}^T = \frac{\sum_{t=1}^T \boldsymbol{\theta}^t}{T}$  and  $\bar{\mathbf{x}}^T = \frac{\sum_{t=1}^T \mathbf{x}^{t+1}}{T}$  to get the near optimal points. It is a standard technique for analyzing stochastic gradient methods [7]. The convergence gap  $\max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\bar{\boldsymbol{\theta}}^T, \mathbf{x}^*) - \max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}^*, \mathbf{x}^*)$  is upper bounded by the regret  $R(T)$

$$R(T) = \sum_{t=1}^T \left[ \max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \phi(\boldsymbol{\theta}^*, \mathbf{x}^*) \right]. \quad (9)$$

**Lemma 4.1.** *For  $\ell$  satisfying assumption 4.1, the objective function  $\phi$  defined in Equation (7)*

$$\max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\bar{\boldsymbol{\theta}}^T, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \max_{\mathbf{x}^* \in \mathcal{B}_{\infty}(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}^*, \mathbf{x}^*) \leq \frac{R(T)}{T}$$

**Adaptive Stochastic Gradient Descent Block Coordinate Ascent (ASGDBCA).** ATAS can be formulated as ASGDBCA, which randomly picks an instance  $\mathbf{x}_k$  at the step  $t$ , applying stochastic gradient descent to the parameter  $\boldsymbol{\theta}$  and adaptive block coordinate ascent to the input  $\mathbf{x}$ . Unlike SGDA [17], where all dimensions of  $\mathbf{x}$  get updated in each iteration, ASGDBCA only updates some dimensions of  $\mathbf{x}$ . ASGDBCA first calculates the pre-conditioner  $v_i^t$  as

$$v_k^{t+1} = \begin{cases} \beta v_i^t + (1 - \beta) \|\nabla_{\mathbf{x}'_i} \ell(\mathbf{x}'_i, y_k; \boldsymbol{\theta}^t)\|_2^2 & i = k \\ v_i^t & i \neq k \end{cases}, \quad \hat{v}_i^{t+1} = \max(\hat{v}_i^t, v_i^{t+1}).$$

Then  $\mathbf{x}$ ,  $\boldsymbol{\theta}$  are optimized with

$$\mathbf{x}_i^{t+1} = \begin{cases} \Pi_{\mathcal{B}_{\infty}(\mathbf{x}_i, \varepsilon)} \left[ \mathbf{x}_i^t + \frac{\eta_x}{\sqrt{\hat{v}_i^{t+1}}} \nabla_{\mathbf{x}'_i} \ell(\mathbf{x}'_i, y_i; \boldsymbol{\theta}^t) \right] & i = k \\ \mathbf{x}_i^t & i \neq k \end{cases}, \quad \boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_{\theta} \nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}_k^{t+1}, y_k; \boldsymbol{\theta}^t).$$

The difference between ASGDBCA and ATAS is  $\hat{v}_k^t$ . To prove the convergence of ASGDBCA, the pre-conditioner needs to be non-decreasing. Otherwise, ATAS may not converge like ADAM [22]. However, the non-convergent version of ADAM actually works better for neural networks in practice [14]. Therefore, ATAS still uses  $v_k^t$  as the pre-conditioner.

**Theorem 4.1** (Regret Bound for ASGDBCA). *Under Assumption 4.1, with  $\eta_{\theta} = \frac{D_{\theta,2}}{G_{\theta,2}\sqrt{T}}$  and  $\eta_x = \frac{\sqrt{d}D_{x,\infty}}{\sqrt{T(1-\beta)^{-1/4}}}$ , the regret of ASGDBCA is bounded by:*

$$R^{ASGDBCA}(T) \leq G_{\theta,2} D_{\theta,2} \sqrt{T} + \frac{D_{x,\infty} \sum_{i=1}^n G_{x_i,2} \sqrt{dT}}{n(1-\beta)^{1/4}} + \frac{dL_x D_{x,\infty}^2}{2n^2 \sqrt{1-\beta}}$$

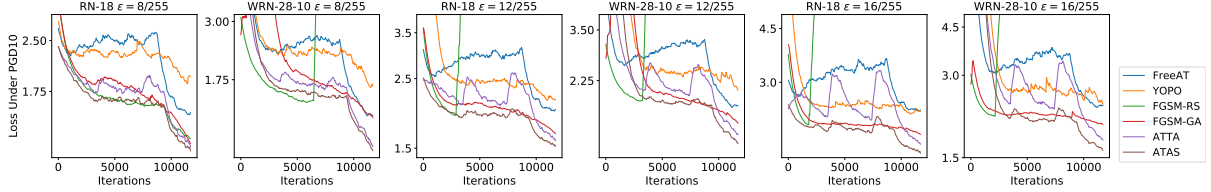


Figure 4: Robust training cross-entropy loss under PGD10 of CIFAR10 with different network architectures and adversarial budgets. The curve is smoothed to clearly show the convergence.

**Comparison with the Non-adaptive Version.** The non-adaptive version of ATAS is ATTA, which can be formulated as the Stochastic Gradient Descent Block Coordinate Ascent (SGDBCA):

$$\mathbf{x}_i^{t+1} = \begin{cases} \Pi_{\mathcal{B}_\infty(\mathbf{x}_i, \varepsilon)}[\mathbf{x}_i^t + \eta_x \nabla_{\mathbf{x}_i^t} \ell(\mathbf{x}_i^t, y_i; \boldsymbol{\theta}^t)] & i = k \\ \mathbf{x}_i^t & i \neq k \end{cases}, \quad \boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_\theta \nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}_k^{t+1}, y_k; \boldsymbol{\theta}^t),$$

**Theorem 4.2** (Regret Bound for SGDBCA). *Under assumption 4.1, with constant learning  $\eta_\theta = \frac{D_{\theta,2}}{G_{\theta,2}\sqrt{T}}$  and  $\eta_x = \frac{\sqrt{nd}D_{x,\infty}}{G_{x,2}\sqrt{T}}$ , the regret  $R^{SGDBCA}(T)$  of SGDBCA is bounded by:*

$$R^{SGDBCA}(T) \leq G_{\theta,2}D_{\theta,2}\sqrt{T} + G_{x,2}D_{x,\infty}\sqrt{\frac{dT}{n}} + \frac{dL_xD_{x,\infty}^2}{2n}$$

Theorem 4.1 and 4.2 shows that ASGDBCA converges faster than SGDBCA. When  $T$  is large, the third term of the regret in both SGDBCA and ASGDBCA is negligible. Consider their first terms are the same, the main difference is the regret bound about  $\mathbf{x}$  in the second term:  $G_{x,2}D_{x,\infty}\sqrt{\frac{dT}{n}}$  and  $\frac{D_{x,\infty}\sum_{i=1}^n G_{x_i,2}\sqrt{dT}}{n(1-\beta)^{1/4}}$ . The ratio between them is

$$\text{Ratio} = \frac{1}{(1-\beta)^{1/4}} \sqrt{\frac{\sum_{i=1}^n G_{x_i,2}^2}{n} / \left(\frac{\sum_{i=1}^n G_{x_i,2}}{n}\right)^2}$$

The Cauchy-Schwarz inequality indicates the ratio is always larger than 1. The gap between ASGDBCA and SGDBCA gets larger when  $G_{x_i,2}$  has long-tailed distribution, which demonstrates the relatively faster convergence of ATAS than the non-adaptive counterparts. We show the empirical histogram of  $G_{x_i,2}$  of a RN-18 and the ratio in Figure 6 in the Appendix, which demonstrates the long-tailed distribution for common datasets.

## 5 Experiments

**Baselines.** We compare ATAS with the SOTA fast AT algorithms including FreeAT [25], YOPO [32], FGSM-RS [30], FGSM-GA [2], SSAT [13] and ATTA [34]. We also compare ATAS with standard AT whose inner maximization is solved by PGD10, providing a reference for the ideal performance.

**Attack Methods.** We consider three attacks: PGD10, PGD50 [19] and AutoAttack (AA) [5]. Square Attack, a black-box attack, is included in AutoAttack to eliminate the effect of gradient masking.

**Experimental Settings.** ATAS uses the techniques proposed in ATTA [34]: the adversarial perturbations are transformed according the data augmentation and get reset every several epochs. And the previous initialization is stored in the GPU memory, brings negligible storing latency to ATAS. We consider adversarial attacks with the  $\ell_\infty$ -norm budget. We evaluate fast AT algorithms on CIFAR10 and CIFAR100 [15] with WideResNet-28-10 (WRN-28-10) [31] and ResNet-18 (RN-18), and on ImageNet [6] with ResNet-18 (RN-18) and ResNet-50 (RN-50). While early stopping is widely used in the standard AT [23], the computational overhead to perform PGD attack on a separate validation set is large. Besides, considering the small budget of training time in fast AT, even if early stopping is applied to terminate the training before catastrophic overfitting occurs, the training is far from convergence, resulting in poor performance [2]. Therefore, we follow the previous works [2, 30, 34] and do not use early stopping. We set  $\beta = 0.5$  and  $\gamma/c = 16/255$ , which is close to the adversarial budget. And we set  $c = 0.01$  for CIFAR10 and CIFAR100 and  $c = 0.1$  for ImageNet. More detailed experiment settings are in Appendix C. Additional experiments are available in the Appendix B.  $\hat{\mathbf{a}}$

Table 1: Accuracy and training time of different methods on CIFAR10, CIFAR100 and ImageNet. ATAS improves the robust accuracy under various attacks including PGD10, PGD50 and AutoAttack (AA). The method “*PGD10*” refers to the standard AT using PGD10 for the inner maximization. Note that, we do not have enough computational resources to perform standard AT and SSAT on ImageNet because of computational complexity. Besides, we are unable to train the ResNet-50 on ImageNet with FGSM-GA as its memory requirement exceeds the maximum GPU memory of our devices (*i.e.* NVIDIA Tesla V100). For CIFAR10 and CIFAR100, the training time is evaluated on a single GPU. And we use two GPUs to train the models for ImageNet. We use default step size from the original papers for the baselines so that catastrophic overfitting seldom happens in these methods.

(a) CIFAR10 with adversarial budget  $\varepsilon = 8/255$ . The accuracy with  $\varepsilon = 12/255$  and  $16/255$  is in the Table 3

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	<i>80.13</i>	<i>50.59</i>	<i>48.94</i>	<i>45.97</i>	<i>1.23</i>	<i>85.00</i>	<i>55.51</i>	<i>53.53</i>	<i>51.27</i>	<i>8.49</i>
FreeAT	78.37	40.90	39.02	36.00	0.33	84.54	46.09	43.80	41.19	2.31
YOPO	74.72	37.51	35.79	33.21	0.28	82.92	44.62	42.14	40.23	1.90
FGSM-RS	<b>83.99</b>	48.99	46.36	42.95	<b>0.22</b>	80.21	0.01	0.00	0.00	1.67
FGSM-GA	80.10	49.14	47.21	43.44	0.57	75.84	45.57	43.28	39.44	3.82
SSAT	88.83	42.31	38.99	37.06	0.61	90.40	44.04	40.40	38.82	3.53
ATTA	82.16	47.47	45.32	42.51	0.30	85.90	51.52	48.94	46.84	1.70
ATAS	81.22	<b>50.03</b>	<b>48.18</b>	<b>45.38</b>	0.30	<b>85.96</b>	<b>53.43</b>	<b>51.03</b>	<b>48.72</b>	<b>1.63</b>

(b) CIFAR100 with adversarial budget  $\varepsilon = 8/255$ . The accuracy with  $\varepsilon = 4/255$  and  $12/255$  is in the Table 5

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	<i>54.08</i>	<i>28.03</i>	<i>27.23</i>	<i>23.04</i>	<i>1.32</i>	<i>60.04</i>	<i>31.70</i>	<i>30.67</i>	<i>27.11</i>	<i>8.53</i>
FreeAT	50.56	19.57	18.58	15.09	0.33	59.38	24.41	23.00	19.60	2.30
YOPO	51.55	20.65	19.17	16.05	0.29	50.35	19.44	18.36	15.43	1.92
FGSM-RS	<b>59.35</b>	26.40	24.29	19.73	<b>0.21</b>	51.83	0.00	0.00	0.00	<b>1.60</b>
FGSM-GA	50.61	24.48	24.07	19.42	0.57	54.29	25.86	24.56	20.74	3.80
SSAT	71.03	9.79	4.80	1.09	0.62	75.01	0.21	0.01	0.00	3.50
ATTA	57.21	25.76	24.90	21.03	0.28	<b>63.04</b>	28.93	27.18	24.42	1.63
ATAS	55.49	<b>27.68</b>	<b>26.60</b>	<b>22.62</b>	0.31	62.34	<b>29.89</b>	<b>28.35</b>	<b>25.03</b>	1.61

(c) ImageNet with adversarial budget  $\varepsilon = 2/255$ . The accuracy with  $\varepsilon = 4/255$  is available in the Table 6

Methods	ResNet-18					ResNet-50				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
FreeAT	58.80	35.56	34.78	31.77	<b>40.01</b>	65.81	44.12	43.34	40.80	<b>108.3</b>
YOPO	47.69	28.50	28.10	25.22	48.22	55.68	33.46	32.19	29.56	111.8
FGSM-RS	55.26	37.33	36.98	33.28	43.46	67.83	46.12	45.56	43.58	115.0
FGSM-GA	37.01	24.15	24.05	19.98	182.7	/	/	/	/	/
ATTA	58.32	39.62	38.32	36.08	45.83	66.62	48.27	47.65	45.00	111.7
ATAS	<b>61.20</b>	<b>40.84</b>	<b>39.86</b>	<b>37.25</b>	45.70	<b>69.10</b>	<b>49.05</b>	<b>48.05</b>	<b>46.01</b>	120.4

Table 2: Ablation study of hyperparameters  $\gamma$  (left) and  $c$  (right) on CIFAR10 and RN-18 under AA.

$\gamma/0.01 * 255$	12	14	16	18	20	$c$	0.005	0.007	0.01	0.02	0.04
$\varepsilon = 8/255$	45.20	45.21	45.38	45.50	45.60	$\varepsilon = 8/255$	45.01	45.28	45.38	45.52	45.48
$\varepsilon = 12/255$	30.84	31.06	30.56	31.21	31.04	$\varepsilon = 12/255$	30.08	30.80	30.56	30.69	30.52
$\varepsilon = 16/255$	21.38	21.23	21.09	21.13	20.94	$\varepsilon = 16/255$	20.36	20.84	21.09	21.07	20.48

**Convergence.** Figure 4 shows the curve of the training loss  $\max_{\mathbf{x}^*=[\mathbf{x}_1^*, \dots, \mathbf{x}_n^*] \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}, \mathbf{x}^*)$  on CIFAR10 with different network architectures and different adversarial budgets, where  $\mathbf{x}^*$  is approximated by PGD10 and the objective function  $\phi$  is approximated by mini-batches of training instances at each step. ATAS achieves smaller robust training loss at the end of training, demonstrating the faster convergence of ATAS than ATTA and other baselines. We also show the relationship between gradient norm distribution



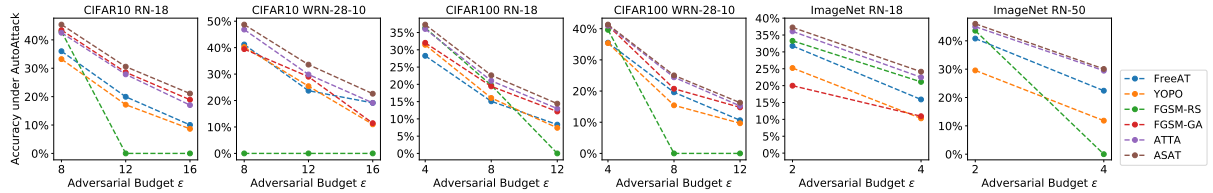


Figure 5: Robust accuracy under AutoAttack for different datasets on different network architectures with varying adversarial budgets. ATAS achieves the highest robust accuracy in these cases. The accuracy numbers can be found in the Appendix B.2.

and convergence gap between ATTA and ATAS in Appendix B.4.

**Robust Accuracy.** We provide our main results in Table 1, showing the robust accuracy of CIFAR10, CIFAR100 and ImageNet, respectively. Figure 5 shows the robust accuracy under AutoAttack for different adversarial budgets, whose numbers are provided in the Appendix B.2.

**CIFAR10 and CIFAR100.** As shown in Table 1a, The robust accuracy of FreeAT and YOPO is much lower than the other methods. While FGSM-RS maintains non-trivial robust accuracy when using RN-18, it suffers from catastrophic overfitting when using large networks such as WRN-28-10. The regularizer in FGSM-GA prevents catastrophic overfitting. However, it may over-regularize the network so that the clean accuracy and the robust accuracy decrease on WRN-28-10. In addition, the regularizer also brings computational overhead: FGSM-GA needs nearly double training time compared with other methods. ATAS achieves the best robust accuracy among all fast AT algorithms while keeping the training time nearly the same. Furthermore, for small networks like RN-18, the performance of ATAS is on par with standard AT (PGD10) but needs only one fifth of the training time. Table 1b shows the robust accuracy on CIFAR100 and ATAS also outperforms other algorithms. Catastrophic overfitting also happens in SSAT even if the losses of inner points are checked.

**ImageNet.** ATTA and ATAS need to memorize the adversarial noise for the whole training set. Since frequently loading and storing from the disks significantly lowers the training speed, all perturbations should be stored in the memory. Thus, we utilize the local property of the adversarial examples [11] and only store the interpolated perturbation in the memory. We resize the perturbations from  $224 \times 224$  to  $32 \times 32$  for storage and up-sample it back when used as the initialization for the next epoch. The detailed algorithm is deferred to the Appendix C.1. Table 1c shows the robust accuracy on ImageNet on  $\varepsilon = 2/255$ . ATAS still has higher robust accuracy than all baselines. FGSM-GA needs to calculate the second order gradient of the parameters, which needs huge amount of GPU memory. Thus, we could not train a big network such as ResNet-50 on ImageNet.

**Robust accuracy at different adversarial budgets.** Figure 5 shows the robust accuracy of fast AT algorithms under AutoAttack on different datasets, network architectures and adversarial budgets. The robust accuracy decreases when enlarging the adversarial budget, but ATAS always outperforms all the baselines for different adversarial budgets, datasets and network architectures. This demonstrates that the improvement of ATAS is consistent.

**Ablation Study.** Table 2 provides the ablation study on hyperparameters, showing that ATAS is not sensitive to them. Besides, as the only different between ATAS and ATTA is the step size, the superior performance of ATAS over ATTA forms an ablation study to demonstrate the effectiveness of the adaptive step size. The changes of gradient norm and step size of ATAS is shown in Appendix B.3.

## 6 Conclusion

In this paper, we investigate catastrophic overfitting from the perspective of training instances and show that instances with large gradient norms are more likely to cause catastrophic overfitting in the single-step fast AT methods. This finding motivates the adaptive training method, ATAS, which applies the adaptive step size of inner maximization inversely proportional to the input gradient norm. We theoretically analyze the convergence of ATAS, showing that our method converges faster than the non-adaptive counterpart especially when the distribution of input gradient norm is long-tailed. Extensive experiments on CIFAR10, CIFAR100 and ImageNet with different network architectures and adversarial budgets show that ATAS mitigates catastrophic overfitting and achieves higher robust accuracy under various strong attacks.

## References

- [1] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*, pages 484–501. Springer, 2020.
- [2] Maksym Andriushchenko and Nicolas Flammarion. Understanding and improving fast adversarial training. *Advances in Neural Information Processing Systems*, 33:16048–16059, 2020.
- [3] Tamer Başar and Geert Jan Olsder. *Dynamic noncooperative game theory*. SIAM, 1998.
- [4] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *International Conference on Machine Learning*, pages 2196–2205. PMLR, 2020.
- [5] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*, pages 2206–2216. PMLR, 2020.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [7] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [8] Biyi Fang and Diego Klabjan. Convergence analyses of online adam algorithm in convex setting and two-layer relu neural network. *arXiv preprint arXiv:1905.09356*, 2019.
- [9] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [10] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten. Countering adversarial images using input transformations. In *International Conference on Learning Representations*, 2018.
- [11] Zhichao Huang, Yaowei Huang, and Tong Zhang. Corrattack: Black-box adversarial attack with structured search. *arXiv preprint arXiv:2010.01250*, 2020.
- [12] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversarial attacks with bandits and priors. In *International Conference on Learning Representations*, 2018.
- [13] Hoki Kim, Woojin Lee, and Jaewook Lee. Understanding catastrophic overfitting in single-step adversarial training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8119–8127, 2021.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [15] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [16] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1778–1787, 2018.
- [17] Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pages 6083–6093. PMLR, 2020.
- [18] Chen Liu, Zhichao Huang, Mathieu Salzmann, Tong Zhang, and Sabine Süsstrunk. On the impact of hard adversarial instances on overfitting in adversarial training. *arXiv preprint arXiv:2112.07324*, 2021.
- [19] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.

- [20] Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics*, pages 1497–1507. PMLR, 2020.
- [21] Geon Yeong Park and Sang Wan Lee. Reliably fast adversarial training via latent adversarial perturbation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7758–7767, 2021.
- [22] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [23] Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In *International Conference on Machine Learning*, pages 8093–8104. PMLR, 2020.
- [24] Tim Roughgarden. Algorithmic game theory. *Communications of the ACM*, 53(7):78–86, 2010.
- [25] Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in Neural Information Processing Systems*, 32, 2019.
- [26] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *International Conference on Learning Representations*, 2018.
- [27] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- [28] T. Tieleman and G. Hinton. Rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 2012.
- [29] 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.
- [30] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In *International Conference on Learning Representations*, 2020.
- [31] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 87.1–87.12. BMVA Press, September 2016.
- [32] Dinghuai Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You only propagate once: Accelerating adversarial training via maximal principle. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pages 7472–7482. PMLR, 2019.
- [34] Haizhong Zheng, Ziqi Zhang, Juncheng Gu, Honglak Lee, and Atul Prakash. Efficient adversarial training with transferable adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1181–1190, 2020.

In the Appendix, we provide additional materials to supplement our main submission. In Appendix A, we provide the proof of Lemma 4.1, Theorem 4.1 and Theorem 4.2 in the main text. In Appendix B, we report additional experimental results including the catastrophic overfitting in ATTA and the robust accuracy of various adversarial budgets. In Appendix C, we provide the detailed hyperparameters and experimental settings.

## A Proof of Section 4.1

We provide the proof of the regret bound in this section.

### A.1 Proof the Lemma 4.1

*Proof.*

$$\begin{aligned}
& \max_{\mathbf{x}^* \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi(\bar{\boldsymbol{\theta}}^T, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \max_{\mathbf{x}^* \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}^*, \mathbf{x}^*) \\
& \leq \max_{\mathbf{x}^* \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi(\bar{\boldsymbol{\theta}}^T, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \phi(\boldsymbol{\theta}^*, \bar{\mathbf{x}}^T) \\
& = \max_{\mathbf{x}^* \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi\left(\frac{\sum_{t=1}^T \boldsymbol{\theta}^t}{T}, \mathbf{x}^*\right) - \min_{\boldsymbol{\theta}^*} \phi\left(\boldsymbol{\theta}^*, \frac{\sum_{t=1}^T \mathbf{x}^{t+1}}{T}\right) \\
& \leq \frac{\min_{\boldsymbol{\theta}^*} \sum_{t=1}^T \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \sum_{t=1}^T \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1})}{T} \\
& \leq \frac{\sum_{t=1}^T (\max_{\mathbf{x}^* \in \mathcal{B}_\infty(\mathbf{x}, \varepsilon)} \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \min_{\boldsymbol{\theta}^*} \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1}))}{T} \\
& = \frac{R(T)}{T}.
\end{aligned} \tag{10}$$

The first and the third inequality follows the optimality condition and the second inequality uses the Jensen inequality.  $\square$

Before moving to the proof of Theorem 4.1 and 4.2, we define several notations of gradients as follows:

$$\begin{aligned}
\hat{g}_\theta^k(\boldsymbol{\theta}, \mathbf{x}) &\equiv \nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}_k, y_k; \boldsymbol{\theta}), \\
g_\theta(\boldsymbol{\theta}, \mathbf{x}) &\equiv \mathbb{E}_k \hat{g}_\theta^k(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \phi(\boldsymbol{\theta}, \mathbf{x}), \\
g_x^k(\boldsymbol{\theta}, \mathbf{x}) &\equiv -\nabla_{\mathbf{x}_k} \ell(\mathbf{x}_k, y_k; \boldsymbol{\theta}), \\
g_x(\boldsymbol{\theta}, \mathbf{x}) &\equiv [g_x^1(\boldsymbol{\theta}, \mathbf{x}), \dots, g_x^n(\boldsymbol{\theta}, \mathbf{x})] = -n \nabla_{\mathbf{x}} \phi(\boldsymbol{\theta}, \mathbf{x}).
\end{aligned} \tag{11}$$

And  $\mathbb{E}_k$  means the average over  $k$ .

### A.2 Proof of Theorem 4.2 (SGDBCA)

*Proof.* Let  $\eta_x = h_x \eta_\theta$ . At step  $t$ , SGDBCA picks a random instance indexed by  $k$  from  $\{1, \dots, n\}$  and updates its adversarial perturbation. Then we have the following inequality:

$$\begin{aligned}
\|\mathbf{x}_k^{t+1} - \mathbf{x}^*\|_2^2 &= \|\Pi_{\mathcal{B}_\infty(\mathbf{x}, \varepsilon)}(\mathbf{x}_k^t - h_x \eta_\theta^t g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)) - \mathbf{x}^*\|_2^2 \\
&\leq \|\mathbf{x}_k^t - h_x \eta_\theta^t g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t) - \mathbf{x}^*\|_2^2.
\end{aligned} \tag{12}$$

Hence

$$\|\mathbf{x}_k^{t+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2^2 + (h_x \eta_\theta)^2 \|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 - 2h_x \eta_\theta g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}_k^t - \mathbf{x}_k^*). \tag{13}$$

Rearranging the inequality, it is easy to get:

$$2\eta_\theta g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}_k^t - \mathbf{x}_k^*) \leq \frac{\|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2^2 - \|\mathbf{x}_k^{t+1} - \mathbf{x}_k^*\|_2^2}{h_x} + h_x (\eta_\theta)^2 \|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2. \tag{14}$$

Similarly, we have:

$$2\eta_\theta g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^t - \boldsymbol{\theta}^*) \leq \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2 + (\eta_\theta)^2 \|\hat{g}_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2. \tag{15}$$

Taking expectation over  $k$  on the left hand side of Equation (14) and Equation (15), we get:

$$\begin{aligned}\mathbb{E}_k [g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}_k^t - \mathbf{x}_k^*)] &= \frac{g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}^t - \mathbf{x}^*)}{n}, \\ \mathbb{E}_k [\hat{g}_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^t - \boldsymbol{\theta}^*)] &= g_\theta(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^t - \boldsymbol{\theta}^*).\end{aligned}\quad (16)$$

Taking expectation over  $k$  on the right hand side of Equation (14), we have:

$$\mathbb{E}_k \left[ \|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2^2 - \|\mathbf{x}_k^{t+1} - \mathbf{x}_k^*\|_2^2 \right] = \frac{1}{n} \left[ \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \right], \quad (17)$$

and

$$\mathbb{E}_k \left[ \|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 \right] = \frac{1}{n} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2. \quad (18)$$

Considering the convex and concave condition of  $\mathbf{x}$  and  $\boldsymbol{\theta}$

$$\begin{aligned}\nabla_{\boldsymbol{\theta}^t} \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^* - \boldsymbol{\theta}^t) &\leq \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1}) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1}), \\ \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^t) + \nabla_{\mathbf{x}^t} \phi(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}^* - \mathbf{x}^t), &\leq \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^t),\end{aligned}\quad (19)$$

we get:

$$\begin{aligned}g_\theta(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^t - \boldsymbol{\theta}^*) + \frac{1}{n} g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}^t - \mathbf{x}^*) \\ \geq \phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1}) + \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1}) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^t).\end{aligned}\quad (20)$$

Combining Eqn (14) to (20), we obtain the following inequality:

$$\begin{aligned}2\eta_\theta (\phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1})) &\leq \mathbb{E}_k \left[ \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2 + (\eta_\theta)^2 \|g_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 + \right. \\ &\quad \left. \frac{1}{nh_x} \left( \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \right) + \right. \\ &\quad \left. \frac{1}{n} h_x (\eta_\theta)^2 \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 + \right. \\ &\quad \left. 2\eta_\theta (\phi(\boldsymbol{\theta}^t, \mathbf{x}^t) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})) \right].\end{aligned}\quad (21)$$

$$\begin{aligned}&\frac{1}{n} h_x (\eta_\theta)^2 \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 + \\ &2\eta_\theta (\phi(\boldsymbol{\theta}^t, \mathbf{x}^t) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})) \left. \right].\end{aligned}\quad (22)$$

Considering the update of  $\mathbf{x}$ , we have:

$$\begin{aligned}&\mathbb{E}_k [\phi(\boldsymbol{\theta}^t, \mathbf{x}^t) - \phi(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})] \\ &\leq \mathbb{E}_k \left[ \frac{L_x}{2n} \|\mathbf{x}_k^t - \mathbf{x}_k^{t+1}\|_2^2 + \frac{1}{n} g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}_k^{t+1} - \mathbf{x}_k^t) \right] \\ &= \mathbb{E}_k \left[ \frac{L_x (h_x \eta_\theta)^2}{2n} \|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 - \frac{h_x \eta_\theta}{n} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 \right] \\ &= \frac{L_x (h_x \eta_\theta)^2}{2n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 - \frac{h_x \eta_\theta}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2.\end{aligned}\quad (23)$$

The above inequality can be rearranged as:

$$\begin{aligned}2\eta_\theta (\phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1})) &\leq \mathbb{E}_k \left[ \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2 + (\eta_\theta)^2 \|g_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 + \right. \\ &\quad \left. \frac{1}{nh_x} \left( \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \right) + \right. \\ &\quad \left. \frac{n-2}{n^2} h_x (\eta_\theta)^2 \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 + \right. \\ &\quad \left. \frac{L_x (h_x)^2 (\eta_\theta)^3}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 \right].\end{aligned}\quad (24)$$

$$\begin{aligned}&\frac{n-2}{n^2} h_x (\eta_\theta)^2 \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 + \\ &\frac{L_x (h_x)^2 (\eta_\theta)^3}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 \left. \right].\end{aligned}\quad (25)$$

Divide both side by  $\eta_\theta$ , then

$$2(\phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1})) \leq \mathbb{E}_k \left[ \frac{1}{\eta_\theta} \left( \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2 \right) + \eta_\theta \|g_\theta(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 + \frac{1}{nh_x\eta_\theta} \left( \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \right) + \frac{n-2}{n^2} h_x \eta_\theta \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 + \frac{L_x(h_x\eta_\theta)^2}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2 \right]. \quad (26)$$

Summing the bound over  $t$ , the regret is bounded by:

$$R(T) \leq \frac{1}{2\eta_\theta} \|\boldsymbol{\theta}^1 - \boldsymbol{\theta}^*\|_2^2 + \frac{1}{2n\eta_x} \|\mathbf{x}^1 - \mathbf{x}^*\|_2^2 + \frac{1}{2} \mathbb{E} \sum_{t=1}^T \eta_\theta \|g_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 + \frac{1}{2} \mathbb{E} \sum_{t=1}^T \frac{[(n-2)\eta_x + L_x\eta_x^2]}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2. \quad (27)$$

With the bound of  $\mathbf{x}$ ,  $\boldsymbol{\theta}$  and their gradients, we can simplify the bound as:

$$R(T) \leq \frac{D_{\theta,2}^2}{2\eta_\theta} + \frac{dD_{x,\infty}^2}{2\eta_x} + \frac{T\eta_\theta G_{\theta,2}^2}{2} + \frac{T\eta_x G_{x,2}^2}{2n} + \frac{TL_x\eta_x^2 G_{x,2}^2}{2n^2}. \quad (28)$$

Using the inequality of arithmetic and geometric means, the optimal choice is  $\eta_\theta = \frac{D_{\theta,2}}{G_{\theta,2}\sqrt{T}}$  and  $\eta_x = \frac{\sqrt{nd}D_{x,\infty}}{G_{x,2}\sqrt{T}}$ , then

$$R(T) \leq G_{\theta,2}D_{\theta,2}\sqrt{T} + G_{x,2}D_{x,\infty}\sqrt{\frac{dT}{n}} + \frac{dL_xD_{x,\infty}^2}{2n}. \quad (29)$$

□

### A.3 Proof of Theorem 4.1 (ASGDBCA)

*Proof.* Let  $\eta_x = h_x\eta_\theta$ . At step  $t$ , ASGDBCA picks a random instance indexed by  $k$  from  $\{1, \dots, n\}$ . Then

$$2g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)^\top (\mathbf{x}_k^t - \mathbf{x}_k^*) \leq \frac{\|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2^2}{\eta_\theta h_x} \sqrt{\hat{v}_k^{t+1}} - \frac{\|\mathbf{x}_k^{t+1} - \mathbf{x}_k^*\|_2^2}{\eta_\theta h_x} \sqrt{\hat{v}_k^{t+1}} + h_x \eta_\theta \frac{\|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2}{\sqrt{\hat{v}_k^{t+1}}}, \quad (30)$$

$$2g_\theta(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})^\top (\boldsymbol{\theta}^t - \boldsymbol{\theta}^*) \leq \frac{\|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2}{\eta_\theta} + \eta_\theta \|g_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2.$$

Let

$$\hat{V}^t = \text{diag}(\underbrace{[\hat{v}_1^t, \dots, \hat{v}_1^t]}_d, \underbrace{[\hat{v}_2^t, \dots, \hat{v}_2^t]}_d, \dots, \underbrace{[\hat{v}_n^t, \dots, \hat{v}_n^t]}_d),$$

denote the pre-conditioner of all the coordinates of  $\mathbf{x}$ . Take the expectation over  $k$  on the right hand side, then

$$\mathbb{E}_k \left[ \frac{\|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2^2 - \|\mathbf{x}_k^{t+1} - \mathbf{x}_k^*\|_2^2}{h_x} \sqrt{\hat{v}_k^{t+1}} \right] = \frac{1}{n} \left[ \|\mathbf{x}^t - \mathbf{x}\|_{\frac{\sqrt{\hat{V}^{t+1}}}{h_x}}^2 - \|\mathbf{x}^{t+1} - \mathbf{x}\|_{\frac{\sqrt{\hat{V}^{t+1}}}{h_x}}^2 \right], \quad (31)$$

and

$$\mathbb{E}_k \left[ \frac{h_x \eta_\theta^2 \|g_x^k(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2^2}{\sqrt{\hat{v}_k^{t+1}}} \right] = \frac{\eta_\theta^2}{n} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{\frac{h_x}{\sqrt{\hat{V}^{t+1}}}}^2. \quad (32)$$

Similar to the proof of SGDBCA, we have:

$$\begin{aligned}
2(\phi(\boldsymbol{\theta}^t, \mathbf{x}^*) - \phi(\boldsymbol{\theta}^*, \mathbf{x}^{t+1})) &\leq \mathbb{E}_k \left[ \frac{1}{\eta_\theta} \left( \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2 - \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|_2^2 \right) + \eta_\theta \|g_\theta(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 + \right. \\
&\quad \frac{1}{nh_x \eta_\theta} \|\mathbf{x}^t - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 + \frac{1}{nh_x \eta_\theta} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 + \\
&\quad \frac{n-2}{n^2} h_x \eta_\theta \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1/2}}^2 + \\
&\quad \left. \frac{L_x (h_x \eta_\theta)^2}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1}}^2 \right]. \tag{33}
\end{aligned}$$

Summing the inequality from 1 to  $T$ , the regret  $R(T) = R_\theta(T) + R_x(T)$  is then upper bounded by:

$$\begin{aligned}
R_\theta(T) &\leq \frac{1}{2\eta_\theta^*} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}\|_2^2 + \frac{\eta_\theta}{2} \mathbb{E} \sum_{t=1}^T \|\hat{g}_\theta^k(\boldsymbol{\theta}^t, \mathbf{x}^{t+1})\|_2^2 \\
&\leq \frac{D_{\theta,2}^2}{2\eta_\theta} + \frac{T\eta_\theta G_{\theta,2}^2}{2}, \tag{34}
\end{aligned}$$

and

$$\begin{aligned}
R_x(T) &\leq \sum_{t=1}^T \left[ \frac{1}{n\eta_x} \left( \|\mathbf{x}^t - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 \right) + \right. \\
&\quad \left. \frac{n-2}{n^2} \eta_x \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1/2}}^2 + \frac{L_x \eta_x^2}{n^2} \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1}}^2 \right]. \tag{35}
\end{aligned}$$

$R_\theta(T)$  is the same as SGD. Using the inequality of arithmetic and geometric means, the optimality is achieved by  $\eta_\theta = \frac{D_{\theta,2}}{G_{\theta,2}\sqrt{T}}$  and we then have:

$$R_\theta(T) \leq G_{\theta,2} D_{\theta,2} \sqrt{T}. \tag{36}$$

For the first term in  $R_x(T)$ , we have:

$$\begin{aligned}
&\sum_{t=1}^T \left( \|\mathbf{x}^t - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 \right) \\
&= \sum_{i=1}^n \left( \sum_{t=2}^T (\sqrt{\hat{v}_i^{t+1}} - \sqrt{\hat{v}_i^t}) \|\mathbf{x}_i^t - \mathbf{x}_i^*\|_2^2 + \sqrt{\hat{v}_i^2} \|\mathbf{x}_i^1 - \mathbf{x}_i^*\|_2^2 \right) \\
&= \sum_{i=1}^n \sum_{j=1}^d \left( \sum_{t=2}^T (\sqrt{\hat{v}_i^{t+1}} - \sqrt{\hat{v}_i^t}) (\mathbf{x}_{i,j}^t - \mathbf{x}_{i,j}^*)^2 + \sqrt{\hat{v}_i^2} (\mathbf{x}_{i,j}^1 - \mathbf{x}_{i,j}^*)^2 \right), \tag{37}
\end{aligned}$$

where  $\mathbf{x}_{i,j}$  means the  $j$ -th coordinate of  $\mathbf{x}_i$ . Since  $\forall i, j, t, |\mathbf{x}_{i,j}^t - \mathbf{x}_{i,j}^*| < D_{x,\infty}$  is assumed, then

$$\begin{aligned}
&\sum_{t=1}^T \left( \|\mathbf{x}^t - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\sqrt{\hat{V}^{t+1}}}^2 \right) \\
&\leq \sum_{i=1}^n \sum_{j=1}^d \left( \sqrt{\hat{v}_i^2} D_{x,\infty}^2 + \sum_{t=2}^T (\sqrt{\hat{v}_i^{t+1}} - \sqrt{\hat{v}_i^t}) D_{x,\infty}^2 \right) \\
&= \sum_{i=1}^n \sum_{j=1}^d \left( \sqrt{\hat{v}_i^{T+1}} D_{x,\infty}^2 \right) \\
&\leq \sum_{i=1}^n \left( d D_{x,\infty}^2 \sqrt{\hat{v}_i^{T+1}} \right) \\
&\leq d D_{x,\infty}^2 \sum_{i=1}^n G_{x_i,2}. \tag{38}
\end{aligned}$$

For the second term of  $R_x(T)$ , we have:

$$\begin{aligned}
& \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1/2}}^2 \\
&= \sum_{i=1}^n \sum_{j=1}^d \frac{(g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2}{\sqrt{v_i^{t+1}}} \\
&\leq \sum_{i=1}^n \sum_{j=1}^d \frac{(g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2}{\sqrt{1-\beta} \sqrt{\sum_{j=1}^d ((g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2)}} \\
&= \frac{1}{\sqrt{1-\beta}} \sum_{i=1}^n \|g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_2 \\
&\leq \frac{1}{\sqrt{1-\beta}} \sum_{i=1}^n G_{x_i,2},
\end{aligned} \tag{39}$$

where  $g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j$  represents the  $j$ -th coordinate of  $g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)$ . Summing over  $t$ , the second term of  $R_x(T)$  is bounded by:

$$\sum_{t=1}^T \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1/2}}^2 \leq \frac{T}{\sqrt{1-\beta}} \sum_{i=1}^n G_{x_i,2}. \tag{40}$$

And the third term of  $R_x(T)$  can be bounded by:

$$\begin{aligned}
& \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1}}^2 \\
&= \sum_{i=1}^n \sum_{j=1}^d \frac{(g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2}{v_i^{t+1}} \\
&\leq \sum_{i=1}^n \sum_{j=1}^d \frac{(g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2}{(1-\beta) \sum_{j=1}^d ((g_x^i(\boldsymbol{\theta}^t, \mathbf{x}^t)_j)^2)} \\
&= \frac{1}{1-\beta}.
\end{aligned} \tag{41}$$

Therefore

$$\sum_{t=1}^T \|g_x(\boldsymbol{\theta}^t, \mathbf{x}^t)\|_{(\hat{V}^{t+1})^{-1}}^2 \leq \frac{T}{1-\beta}. \tag{42}$$

Combining these inequalities,  $R_x(T)$  is bounded by:

$$R_x(T) \leq \frac{T\eta_x}{2n\sqrt{1-\beta}} \sum_{i=1}^n G_{x_i,2} + \frac{D_{x,\infty}^2 d}{2n\eta_x} \sum_{i=1}^n G_{x_i,2} + \frac{L_x \eta_x^2 T}{2n^2(1-\beta)}. \tag{43}$$

Using the inequality of arithmetic and geometric means, the bound on the achieves the minimum when  $\eta_x = \frac{\sqrt{d}D_{x,\infty}(1-\beta)^{1/4}}{\sqrt{T}}$  and

$$R_x(T) \leq \frac{D_{x,\infty} \sum_{i=1}^n G_{x_i,2} \sqrt{dT}}{n(1-\beta)^{1/4}} + \frac{dL_x D_{x,\infty}^2}{2n^2 \sqrt{1-\beta}}. \tag{44}$$

Combining  $R_\theta(T)$  and  $R_x(T)$ , the regret of ASGDBCA is bounded by

$$\begin{aligned}
R(T) &\leq G_{\theta,2} D_{\theta,2} \sqrt{T} + \frac{D_{x,\infty} \sum_{i=1}^n G_{x_i,2} \sqrt{dT}}{n(1-\beta)^{1/4}} + \\
&\quad \frac{dL_x D_{x,\infty}^2}{2n^2 \sqrt{1-\beta}}.
\end{aligned} \tag{45}$$

□



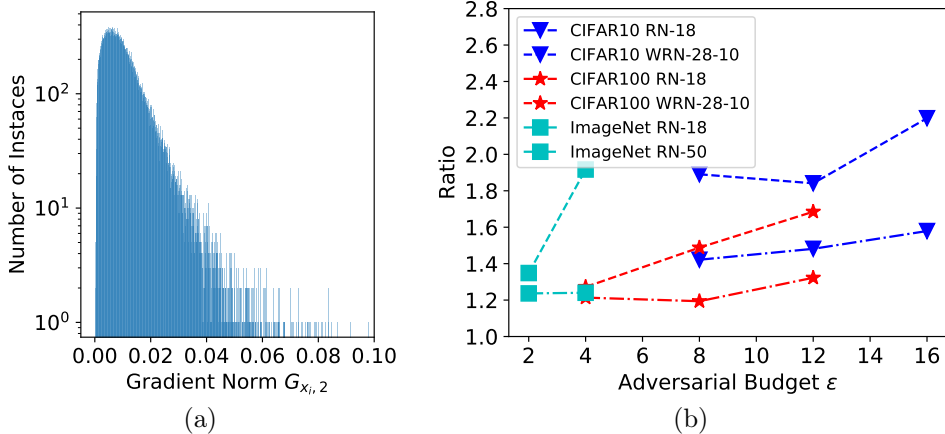


Figure 6: (a) Histogram of  $G_{x_i,2}$  of a ResNet-18 trained CIFAR10 with  $\epsilon = 8/255$ . (b) Ratio  $\frac{\sum_{i=1}^n G_{x_i,2}^2}{(\frac{\sum_{i=1}^n G_{x_i,2}}{n})^2}$  for different datasets and network architectures with different  $\epsilon$ .

## B Additional Experiments

### B.1 Catastrophic Overfitting in ATTA

As shown in Figure 7a, when increasing the step size in ATTA, the loss gap between the ATTA and PGD10 becomes smaller. Furthermore, the robust accuracy also increases when the step size is not overly large in Figure 7b. It shows that large step size also strengthen the attack in ATTA. However, large step also leads to catastrophic overfitting in ATTA. When the step size is overly large in Figure 7b, the robust accuracy against PGD decreases to nearly 0%, indicating catastrophic overfitting.

As we show in the main text, instances with large gradient norm are more likely to cause catastrophic overfitting in ATTA. To verify it, we train a RN-18 with  $\epsilon = 8/255$  on CIFAR10 with ATTA. We record the average gradient norm  $GN(\mathbf{x}_i)$  and divide the subsets according to  $\text{rank}(\mathbf{x}_i)$ :

$$\mathcal{D}_i^j = \{\mathbf{x}_k \mid \frac{10(i-1)}{n} \leq \text{rank}(\mathbf{x}_k) < \frac{10j}{n}\}.$$

The training curves for different subsets are shown in Figure 9. It is the ATTA version of Figure 2 in the main text. When training with  $\epsilon = 12/255$   $\alpha = 15/255$  or  $\epsilon = 12/255$   $\alpha = 18/255$ , the robust accuracy of the subsets of large gradient norm ( $\mathcal{D}_7^{10}, \mathcal{D}_8^{10}, \mathcal{D}_9^{10}$ ) suddenly decreases to nearly 0%, indicating the phenomenon of catastrophic overfitting. By contrast, catastrophic overfitting does not occur when training with the subsets of small gradient norm ( $\mathcal{D}_1^2, \mathcal{D}_1^3, \mathcal{D}_1^4$ ). The observation is the same as FGSM-RS in the main text.

Adaptive step sizes in ATAS allow larger step sizes without causing catastrophic overfitting. In Figure 8, we show the comparison between ATTA and ATAS. Even if ATAS has larger step size than ATTA, it does not suffer from catastrophic overfitting like ATTA.

### B.2 Robust Accuracy of Various Adversarial Budget

The robust accuracy for CIFAR10 with  $\epsilon = 8/255, 12/255, 16/255$ , CIFAR100 with  $\epsilon = 4/255, 8/255, 12/255$  and ImageNet with  $\epsilon = 2/255, 4/255$  are provided in Table 3, Table 5 and Table 6 respectively. ATAS achieves the best robust accuracy in all these experiments with different datasets, network architectures and adversarial budgets.

### B.3 Steps size and Gradient Norm

Figure 10 plots the changes of gradient norm and step size for ATAS after warm up. We divide CIFAR10 into 10 subsets according to their gradient norm and plot the gradient norm and step size for  $\mathcal{D}_1^1, \mathcal{D}_5^5$  and  $\mathcal{D}_{10}^{10}$ , which has the smallest, medium and largest input gradient norm among 10 subsets. The figure shows that the input gradient norm and step size is relatively stable for each subset along the training process. It shows that the input gradient norm is more like a property of training instances themselves,

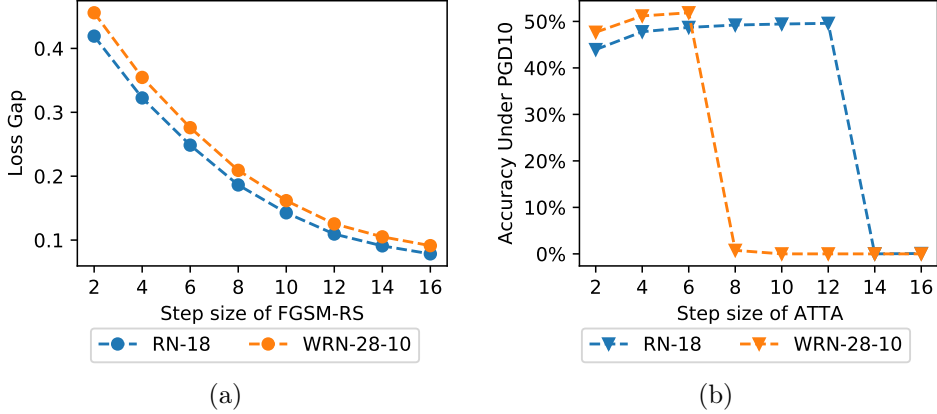


Figure 7: (a) The loss gap of training instances between PGD10 and ATTA attack  $\ell(\mathbf{x}^{\text{PGD}}, y) - \ell(\mathbf{x}^{\text{FGSM-RS}}, y)$  with different step sizes for a ATTA trained robust model; (b) The test robust accuracy of the models trained by ATTA with different step sizes.

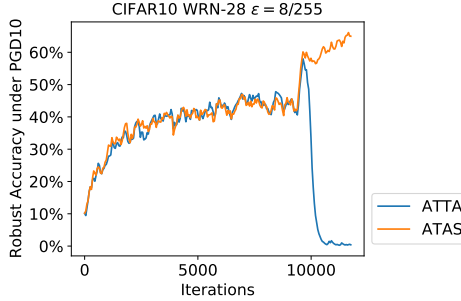


Figure 8: Training accuracy under PGD10 of ATTA ( $\alpha = 8/255$ ) and ATAS. Even if the average step size of ATAS ( $\bar{\alpha} = 9.3/255$ ) is larger than ATTA ( $\alpha = 8/255$ ), catastrophic overfitting does not occur in ATAS.

which is consistent with our motivation. It is worth noting that the sudden changes of gradient norm is the result of initialization reset used in ATTA.

## B.4 Convergence Gap

In Table 4, we show the relationship between the *Ratio*

$$\text{Ratio} = \frac{1}{(1 - \beta)^{\frac{1}{4}}} \sqrt{\frac{\sum_{i=1}^n G_{x_i,2}^2}{n} / \left(\frac{\sum_{i=1}^n G_{x_i,2}}{n}\right)^2}$$

and the convergence gap  $\ell_{\text{ATTA}} - \ell_{\text{ATAS}}$  and convergence ratio  $\ell_{\text{ATTA}}/\ell_{\text{ATAS}}$  in the last epoch of training. Here,  $\ell$  is the loss of each method. The ratio is obtained from Figure 6b for CIFAR10 with ResNet-18. It shows that larger *Ratio* (more long-tailed distribution) leads to larger convergence gap between ATTA and ATAS.

## C Details about the Experiments

### C.1 Algorithms for ATAS in ImageNet

In the experiments of ATTA and ATAS, we utilize the local property of the adversarial examples [11, 12] and only store the interpolated perturbation in the memory. We resize the perturbations from  $224 \times 224$  to  $32 \times 32$  for storage in the memory and up-sample it back when using it as the initialization for the next epoch. The detailed algorithm is shown in Algorithm 2.

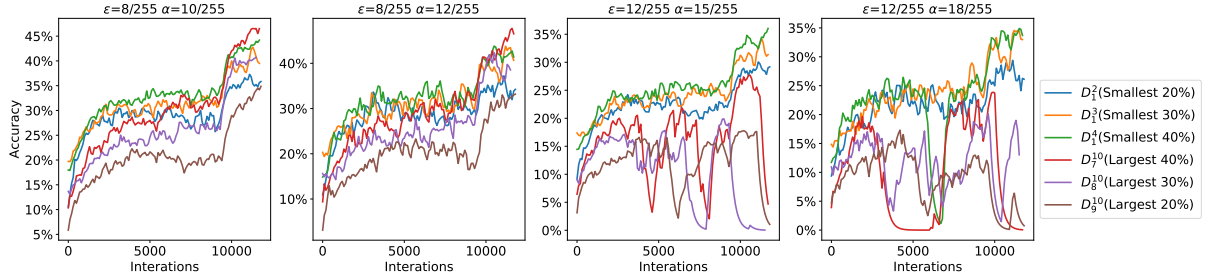


Figure 9: The robust training accuracy curve of FGSM-RS trained on different subsets of CIFAR10. The attack is PGD10 and the network is ResNet-18. The adversarial budgets and the step sizes are shown on top of each figure.

Table 3: Robust accuracy and training time of different fast AT methods on CIFAR10 with  $\varepsilon = 8/255, 12/255, 16/255$ .

(a) CIFAR10 with adversarial budget  $\varepsilon = 8/255$ .

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	80.13	50.59	48.94	45.97	1.23	85.00	55.51	53.53	51.27	8.49
FreeAT	78.37	40.90	39.02	36.00	0.33	84.54	46.09	43.80	41.19	2.31
YOPO	74.72	37.51	35.79	33.21	0.28	82.92	44.62	42.14	40.23	1.90
FGSM-RS	<b>83.99</b>	48.99	46.36	42.95	<b>0.22</b>	80.21	0.01	0.00	0.00	1.67
FGSM-GA	80.10	49.14	47.21	43.44	0.57	75.84	45.57	43.28	39.44	3.82
SSAT	88.83	42.31	38.99	37.06	0.61	90.40	44.04	40.40	38.82	3.53
ATTA	82.16	47.47	45.32	42.51	0.30	85.90	51.52	48.94	46.84	1.70
ATAS	81.22	<b>50.03</b>	<b>48.18</b>	<b>45.38</b>	0.30	<b>85.96</b>	<b>53.43</b>	<b>51.03</b>	<b>48.72</b>	<b>1.63</b>

(b) CIFAR10 with adversarial budget  $\varepsilon = 12/255$

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	70.46	39.39	37.22	32.99	1.23	76.49	43.67	40.93	36.95	8.48
FreeAT	72.92	25.88	22.82	19.95	0.33	79.71	26.31	23.72	18.98	2.33
YOPO	64.21	23.82	22.27	17.16	0.29	75.29	32.27	28.41	25.42	1.92
FGSM-RS	<b>80.78</b>	0.00	0.00	0.00	<b>0.22</b>	79.41	0.00	0.00	0.00	1.66
FGSM-GA	68.62	36.76	33.96	28.57	0.57	72.87	37.98	35.18	29.01	3.82
SSAT	89.08	6.50	1.16	0.03	0.60	91.45	0.07	0.00	0.00	3.50
ATTA	74.46	35.85	31.69	27.85	0.28	<b>80.05</b>	38.29	34.01	29.85	1.63
ATAS	72.58	<b>38.10</b>	<b>35.58</b>	<b>30.56</b>	0.29	78.16	<b>41.88</b>	<b>38.94</b>	<b>33.58</b>	<b>1.62</b>

(c) CIFAR10 with adversarial budget  $\varepsilon = 16/255$

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	61.08	31.42	29.37	23.34	1.32	66.57	35.91	32.89	27.24	8.65
FreeAT	61.05	15.86	12.49	10.04	0.33	67.89	19.07	14.76	12.53	2.34
YOPO	67.45	14.87	12.00	8.66	0.29	62.75	18.30	16.27	11.00	1.91
FGSM-RS	67.75	0.00	0.00	0.00	<b>0.22</b>	60.21	0.00	0.00	0.00	1.67
FGSM-GA	54.07	27.05	25.10	18.92	0.57	16.08	13.38	13.29	11.44	3.77
SSAT	90.41	0.43	0.03	0.00	0.59	91.42	0.00	0.00	0.00	3.52
ATTA	63.37	26.66	23.45	17.02	0.29	<b>72.90</b>	30.11	23.92	19.11	<b>1.65</b>
ATAS	<b>64.11</b>	<b>31.39</b>	<b>28.15</b>	<b>21.09</b>	0.30	70.33	<b>34.32</b>	<b>30.53</b>	<b>22.58</b>	1.68

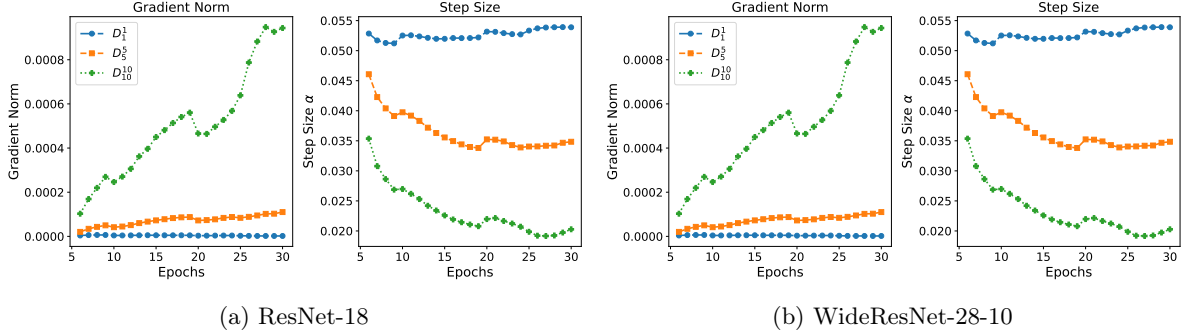


Figure 10: The input gradient norm and step size for CIFAR10 with  $\varepsilon = 8/255$  after warm up.

Table 4: Convergence gap and the ratio on CIFAR10 with ResNet-18.

Ratio	1.4 ( $\varepsilon=8/255$ )	1.5 ( $\varepsilon=12/255$ )	1.6 ( $\varepsilon=16/255$ )
Convergence Gap $\ell_{\text{ATTA}} - \ell_{\text{ATAS}}$	0.05	0.10	0.12
Convergence Ratio $\ell_{\text{ATTA}}/\ell_{\text{ATAS}}$	1.03	1.11	1.13

## C.2 Detailed Hyperparameters for the Experiments

As we focus on fast AT, we reduce the training epochs like [2, 30]. For single-step methods FGSM-RS, FGSM-GA, ATTA and ATAS, the training lasts for 30 epochs on CIFAR10 and CIFAR100, and 90 epochs on ImageNet. For FreeAT and YOPO, we keep the number of the forward-backward passes the same as the single-step methods so that the total training time of these methods will be similar. We use two kinds of learning rate scheduler: piece-wise decay used in [34] and cyclic learning rate used in [30], and choose the best scheduler for each method.

**FreeAT.** We use the default hyperparameters from [25] except training epochs to make fair comparison between different methods. We select the best number of batch replaying from [25]. For CIFAR10 and CIFAR100, we use Free-8 in their paper (Free- $m$  means the number of batch replaying is  $m$ ) and train the network for 10 epochs. For ImageNet, we use Free-4 and train the network for 45 epochs.

**YOPO.** We use YOPO-5-3 in [32] as it achieves the best performance. The training lasts for 12 epochs for CIFAR10 and CIFAR100. For ImageNet, the training lasts for 36 epochs to make the training time similar to other methods. Other hyperparameters are the same as the original paper [32].

**FGSM-RS.** We directly download the code from the official repository [https://github.com/locuslab/fast\\_adversarial](https://github.com/locuslab/fast_adversarial). The training lasts for 30 epochs for CIFAR10 and CIFAR100, and 90 epochs for ImageNet. Following the hyperparameters in the paper, the step size  $\alpha = 1.25\varepsilon$ . Other hyperparameters are the same as their paper.

**FGSM-GA.** We directly download the code from the official repository <https://github.com/tml-epfl/understanding-fast-adv-training>. The training lasts for 30 epochs for CIFAR10 and CIFAR100, and 90 epochs for ImageNet. Other hyperparameters are the same. For the experiments not involved in their paper, we keep them same as the experiments of CIFAR10 except for the hyperparameter  $\lambda$  balancing the gradient align regularizer, which also varies for different datasets and adversarial budgets in their code.  $\lambda$  for CIFAR10 is provided in their code. For CIFAR100 and ImageNet, we run several experiments and provide the result with best  $\lambda$ . These  $\lambda$  are provided in Table 7.

**SSAT.** We directly download the code from the official repository <https://github.com/Harry24k/catastrophic-overfitting>. The training lasts for 30 epochs for CIFAR10 and CIFAR100. And we use the check points  $c = 3$ , which achieves the best performance in their paper.

**ATTA.** We follow the hyperparameters setting for ATTA-1 in [34] and set the step size  $\alpha = 4/255$ . We reduce the number of training epochs to 30 for CIFAR10 and CIFAR100. And the epochs of piece-wise learning rate are rescheduled accordingly. The learning rate  $\eta$  starts at 0.1 and decays to 0.01 and 0.001 at the 24th and 28th epochs. The training of ImageNet lasts for 90 epochs and the learning rate also starts at 0.1 and decays to 0.01 and 0.001 at the 50th and 75th epochs. The weight decay is  $5 \times 10^{-4}$  for CIFAR10 and CIFAR100. For ImageNet, it is  $1 \times 10^{-4}$ . The batch size is 128 for all the experiments. Other hyperparameters are the same as their paper.

**ATAS.** The hyperparameters  $\gamma$  and  $c$  are used to control the minimum and maximum step size for the training instances. When the moving average of gradient norm  $v_i^j \rightarrow 0$ , the step size  $\alpha_i^j = \gamma/c$ . We choose

---

**Algorithm 2** ATAS for ImageNet

---

**Input:** Training set  $\mathcal{D}$ , The model  $f_{\theta}$  with loss function  $\ell$ , Adversarial budget  $\varepsilon$ , Hyperparameters  $\gamma, \eta, c, N$

**Output:** Optimized model  $f_{\theta^*}$

- 1:  $v_i^0 = 0$  for  $i = 1, \dots, n$
  - 2:  $\delta_i^0 = \text{Uniform}(-\varepsilon, \varepsilon)$  for  $i = 1, \dots, n$
  - 3: Resize  $\delta_i^0$  to  $32 \times 32$  for  $i = 1, \dots, n$  and store them.
  - 4: **for**  $j = 1$  to  $N$  **do**
  - 5:   **for**  $\mathbf{x}_i, y_i \in \mathcal{D}$  **do**
  - 6:     Resize  $\delta_i^{j-1}$  to  $224 \times 224$
  - 7:      $\mathbf{x}_i^{j-1} = \mathbf{x}_i + \delta_i^{j-1}$
  - 8:      $v_i^j = \beta v_i^{j-1} + (1 - \beta) \|\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y_i; \theta)\|_2^2$
  - 9:      $\alpha_i^j = \gamma / (c + \sqrt{v_i^j})$
  - 10:      $\mathbf{x}_i^j = \Pi_{\mathcal{B}_p(\mathbf{x}_i, \varepsilon)}[\mathbf{x}_i^{j-1} + \alpha_i^j \cdot \text{sgn}(\nabla_{\mathbf{x}_i^{j-1}} \ell(\mathbf{x}_i^{j-1}, y; \theta))]$
  - 11:      $\theta = \theta - \eta \nabla_{\theta} \ell(\mathbf{x}_i^j, y; \theta)$
  - 12:      $\delta_i^j = \mathbf{x}_i^j - \mathbf{x}_i$
  - 13:     Resize  $\delta_i^j$  to  $32 \times 32$  and store it.
  - 14:   **end for**
  - 15: **end for**
- 

$\gamma/c = 16/255$ , which is close to the adversarial budget. And  $c$  should be close to the magnitude of  $v_i^j$ . As the gradient norm increases with the dimension of the inputs,  $c$  should be larger for ImageNet. Therefore, we set  $c = 0.01$  and for CIFAR10 and CIFAR100 and we let  $c = 0.1$  for ImageNet. Momentum of gradient norm  $\beta$  is set to 0.5 for all the experiments. ATAS is not sensitive to the choice of hyperparameters. Other hyperparameters are the same as ATTA.

### C.3 Environments of the Experiments

All the training time is evaluated on a machine with *Intel Xeon 8255C* and *NVIDIA Tesla V100*. For CIFAR10 and CIFAR100, we use a single GPU. For ImageNet, we use two GPUs. We run all the experiments with *Pytorch 1.4*.

Table 5: Robust accuracy and training time of different fast AT methods on CIFAR100 with  $\varepsilon = 4/255, 8/255, 12/255$ .

(a) CIFAR100 with adversarial budget  $\varepsilon = 4/255$

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	<i>63.22</i>	<i>41.18</i>	<i>40.57</i>	<i>37.75</i>	<i>1.22</i>	<i>69.03</i>	<i>45.27</i>	<i>44.44</i>	<i>43.30</i>	<i>8.61</i>
FreeAT	54.38	32.21	31.68	28.26	0.32	63.91	39.39	38.64	35.52	2.29
YOPO	59.04	34.55	34.02	31.45	0.28	64.60	38.92	38.27	35.39	1.89
FGSM-RS	<b>65.50</b>	39.41	38.50	36.35	<b>0.22</b>	69.62	42.18	41.35	39.64	<b>1.60</b>
FGSM-GA	57.33	35.49	35.01	32.03	0.57	68.45	<b>51.92</b>	<b>44.09</b>	41.03	3.80
SSAT	70.81	33.17	31.09	29.81	0.60	74.43	36.51	34.55	33.34	3.52
ATTA	64.28	39.55	39.20	36.03	0.29	<b>69.51</b>	<b>44.36</b>	42.66	40.99	1.66
ATAS	63.79	<b>40.68</b>	<b>40.02</b>	<b>37.30</b>	0.30	69.64	44.34	43.36	<b>41.32</b>	1.66

(b) CIFAR100 with adversarial budget  $\varepsilon = 8/255$

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	<i>54.08</i>	<i>28.03</i>	<i>27.23</i>	<i>23.04</i>	<i>1.32</i>	<i>60.04</i>	<i>31.70</i>	<i>30.67</i>	<i>27.11</i>	<i>8.53</i>
FreeAT	50.56	19.57	18.58	15.09	0.33	59.38	24.41	23.00	19.60	2.30
YOPO	51.55	20.65	19.17	16.05	0.29	50.35	19.44	18.36	15.43	1.92
FGSM-RS	<b>59.35</b>	26.40	24.29	19.73	<b>0.21</b>	51.83	0.00	0.00	0.00	<b>1.60</b>
FGSM-GA	50.61	24.48	24.07	19.42	0.57	54.29	25.86	24.56	20.74	3.80
SSAT	71.03	9.79	4.80	1.09	0.62	75.01	0.21	0.01	0.00	3.50
ATTA	57.21	25.76	24.90	21.03	0.28	<b>63.04</b>	28.93	27.18	24.42	1.63
ATAS	55.49	<b>27.68</b>	<b>26.60</b>	<b>22.62</b>	0.31	62.34	<b>29.89</b>	<b>28.35</b>	<b>25.03</b>	1.61

(c) CIFAR100 with adversarial budget  $\varepsilon = 12/255$

Methods	ResNet-18					WideResNet-28-10				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
<i>PGD10</i>	<i>44.31</i>	<i>20.41</i>	<i>19.22</i>	<i>15.41</i>	<i>1.34</i>	<i>50.30</i>	<i>23.81</i>	<i>22.55</i>	<i>18.13</i>	<i>8.56</i>
FreeAT	41.05	11.85	10.67	8.33	0.32	46.54	14.95	13.07	10.64	2.30
YOPO	44.69	10.52	9.20	7.41	0.29	54.13	13.19	11.76	9.68	1.92
FGSM-RS	32.78	0.00	0.00	0.00	<b>0.22</b>	38.74	0.00	0.00	0.00	<b>1.60</b>
FGSM-GA	39.77	17.06	16.07	12.14	0.57	51.05	20.54	19.37	14.77	3.80
SSAT	71.38	3.00	1.18	0.09	0.60	75.50	0.00	0.00	0.00	3.56
ATTA	<b>50.55</b>	18.58	16.59	12.97	0.28	<b>56.46</b>	20.82	18.17	15.24	1.63
ATAS	47.14	<b>19.73</b>	<b>18.39</b>	<b>14.41</b>	0.31	53.70	<b>22.53</b>	<b>20.95</b>	<b>16.27</b>	1.63

Table 6: Robust accuracy and training time of different fast AT methods on ImageNet with  $\varepsilon = 2/255, 4/255$ .

(a) ImageNet with adversarial budget  $\varepsilon = 2/255$

Methods	ResNet-18					ResNet-50				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
FreeAT	58.80	35.56	34.78	31.77	<b>40.01</b>	65.81	44.12	43.34	40.80	<b>108.3</b>
YOPO	47.69	28.50	28.10	25.22	48.22	55.68	33.46	32.19	29.56	111.8
FGSM-RS	55.26	37.33	36.98	33.28	43.46	67.83	46.12	45.56	43.58	115.0
FGSM-GA	37.01	24.15	24.05	19.98	182.7	/	/	/	/	/
ATTA	58.32	39.62	38.32	36.08	45.83	66.62	48.27	47.65	45.00	111.7
ATAS	<b>61.20</b>	<b>40.84</b>	<b>39.86</b>	<b>37.25</b>	45.70	<b>69.10</b>	<b>49.05</b>	<b>48.05</b>	<b>46.01</b>	120.4

(b) ImageNet with adversarial budget  $\varepsilon = 4/255$

Methods	ResNet-18					ResNet50				
	Clean	PGD10	PGD50	AA	Time(h)	Clean	PGD10	PGD50	AA	Time(h)
FreeAT	56.99	20.75	18.86	15.90	40.18	64.25	27.95	25.46	22.40	109.6
YOPO	33.72	13.36	13.01	10.30	48.33	37.62	14.77	13.37	11.83	111.9
FGSM-RS	49.73	26.48	25.70	21.11	<b>38.84</b>	66.75	1.08	0.13	0.00	<b>103.6</b>
FGSM-GA	29.34	15.58	15.42	10.94	180.2	/	/	/	/	/
ATTA	54.25	27.31	26.97	22.47	47.95	63.28	35.13	33.37	29.46	114.8
ATAS	<b>55.69</b>	<b>29.23</b>	<b>28.13</b>	<b>24.13</b>	44.15	<b>65.26</b>	<b>35.74</b>	<b>33.58</b>	<b>30.07</b>	110.7

Table 7: Hyperparameter  $\lambda$  for FGSM-GA

(a) CIFAR100				(b) ImageNet		
$\varepsilon$	4/255	8/255	12/255	$\varepsilon$	2/255	4/255
$\lambda$	0.2	0.5	1.0	$\lambda$	0.005	0.01