
A²: Efficient Automated Attacker for Boosting Adversarial Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Based on the significant improvement of model robustness by AT (Adversarial
2 Training), various variants have been proposed to further boost the performance.
3 Well-recognized methods have focused on different components of AT (e.g., de-
4 signing loss functions and leveraging additional unlabeled data). It is generally
5 accepted that stronger perturbations yield more robust models. However, how to
6 generate stronger perturbations efficiently is still missed. In this paper, we propose
7 an efficient automated attacker called A² to boost AT by generating the optimal
8 perturbations on-the-fly during training. A² is a parameterized automated attacker
9 to search in the attacker space for the best attacker against the defense model and
10 examples. Extensive experiments across different datasets demonstrate that A²
11 generates stronger perturbations with low extra cost and reliably improves the
12 robustness of various AT methods against different attacks.

13 1 Introduction

14 DNNs (Deep Neural Networks) are extremely vulnerable to imperceptible perturbations despite their
15 success in a wide variety of applications [He et al., 2016, Kenton and Toutanova, 2019, Guo et al.,
16 2017]. In particular, adding small but carefully chosen deviations to the input, called adversarial
17 perturbations, can cause DNNs to make incorrect predictions with high confidence. It indicates that
18 models trained by minimizing the empirical risk are not intrinsically robust. To explicitly improve
19 robustness, AT (Adversarial Training), where a defense model is trained on worst-case adversarial
20 perturbations generated by an attacker, was developed and proved to be effective.

21 Based on the significant improvement of models’ robustness, various methods have been pro-
22 posed, which focus on different components of AT: analyzing the robustness of neural architectures
23 Zagoruyko and Komodakis [2016], designing loss functions such as TRADES [Zhang et al., 2019]
24 and MART [Wang et al., 2019], perturbing the model to regularize the loss landscape’s flatness (i.e.,
25 AWP Wu et al. [2020]) and leveraging unlabeled data (i.e., RST [Carmon et al., 2019]). Goyal et al.
26 [2020] compares the performance of each combination of most components and achieves the SOTA
27 (state-of-the-art) performance.

28 All the above AT methods use PGD^K [Madry et al., 2018], which is a K -step stack of FGSM [Good-
29 fellow et al., 2015], as the attacker to generate perturbations for each example against the defense
30 model. As the key of AT, stronger perturbations yield more robust models. However, there is a
31 trade-off between the strength of the perturbation and the training efficiency. Increasing the attack
32 step K strengthens the perturbations, but linearly increases the training overhead [Goyal et al., 2020].
33 Likewise, the huge overhead prevents the use of SOTA adversarial attackers [Croce and Hein, 2020,
34 Yao et al., 2021]. To achieve a balance of robustness and efficiency, manually tuning the attacker
35 (e.g., step size and attack method in each step) is of great concern. R+FGSM [Wong et al., 2019]

36 first randomly initializes a small perturbation, and then applies FGSM with the tuned step size.
 37 Surprisingly, AT with R+FGSM is as effective as PGD^K but has a significantly lower cost.

38 However, given the novel dataset, tuning the attacker manually is a challenging task requiring expert
 39 knowledge. Moreover, the best attacker is fine-grained to each example and current model during
 40 adversarial training. Manual coarse-grained tuning for the whole training (e.g., fixed attack method
 41 and step size for all examples) is sub-optimal and prevents further improvement of robustness.

42 Inspired by AutoML (Automated Machine Learning [Zoph and Le, 2017, Liu et al., 2018]), we
 43 propose an efficient automated attacker called A² to boost AT. A² is a parameterized attacker, which
 44 can automatically tune itself on-the-fly during training to generate worst-case perturbations for each
 45 example. First, we design a general attacker space by referring to existing attackers. The attacker
 46 space is stacked by one-step attacker cells. Each cell consists of a perturbation block and a step
 47 size block. Then, we employ a parameterized attacker to search for operations in each block and
 48 construct the attacker for each example that maximizes the model loss. Specifically, we leverage
 49 the attention mechanism to calculate the score of each operation. For continuous operations, we
 50 sum up the operations using the normalized scores as weights. For discrete operations, we use the
 51 reparameterization trick [Jang et al., 2017] to sample an operation from the corresponding block. In
 52 this way, the constructed attacker generates worst-case adversarial perturbations to train the defense
 53 model. Meanwhile, A² is differentiable and can be optimized with respect to the model loss by
 54 gradient descent.

55 We conduct extensive experiments to verify the effectiveness and efficiency of A². Compared
 56 with PGD, A² can find better perturbations for different models trained with various AT methods.
 57 The results demonstrate that 20-step A² generates better perturbations than PGD¹⁰⁰. Moreover,
 58 we combine A² with other AT variants, and the robustness of models with different architectures
 59 on various datasets is generally improved under strong attacks (e.g., classical C&W and SOTA
 60 AutoAttack). We also show that A² is insensitive to its hyperparameters.

61 To summarize, our main contributions can be highlighted as follows:

- 62 • We propose an efficient automated attacker called A², which can generate worst-case
 63 perturbations on-the-fly during training to improve robustness.
- 64 • In A², we design an attacker space by summarizing the existing attackers and employ a
 65 differentiable method to construct the most adversarial attacker for each example according
 66 to the attention mechanism.
- 67 • Extensive experimental results across different datasets and neural architectures demonstrate
 68 that A² improves the model’s robustness by generating stronger perturbations in the *inner*
 69 *maximization*. Moreover, A² can be flexibly combined with different AT methods, showing
 70 good generality.

71 2 Preliminary: Adversarial Training

72 Let $D = (\mathbf{X}, Y) = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be a dataset with $\mathbf{x}_i \in \mathbb{R}^d$ as a natural example and $y_i \in \{1, \dots, C\}$
 73 as its associated label. We measure the performance of a DNN classifier f parametrized with θ using
 74 a suitable loss function l , denoted as $\mathbb{E}_{(\mathbf{x}_i, y_i) \in D} [l(f_\theta(\mathbf{x}_i), y_i)]$. AT [Madry et al., 2018] formulates a
 75 saddle point problem whose goal is to find the model parameters θ that minimize the adversarial risk
 76 in the *outer minimization* (the example’s index i is omitted for brevity):

$$\underbrace{\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \in D} \left[\overbrace{\max_{\delta \in \mathbb{S}} l(f_\theta(\mathbf{x} + \delta), y)}^{\text{inner maximization}} \right]}_{\text{outer minimization}} \quad (1)$$

77 where \mathbb{S} defines the set of allowed perturbations. The perturbation is usually constrained by L_p norm
 78 with a bound ϵ , i.e. $\mathbb{S} = \{\delta \mid \|\delta\|_p \leq \epsilon\}$.

79 The *inner maximization* aims to find an adversarial perturbation against the example that achieves a
 80 high loss for the defense model. However, it is NP-hard to find the optimum of the *inner maximization*.
 81 Various gradient-based attackers have been proposed to approximate its solution, and we classify

82 them according to the number of steps in gradient ascent. One-step attackers [Goodfellow et al., 2015,
83 Miyato et al., 2017] generate adversarial perturbations as:

$$\delta^* \approx \Pi_{\mathbb{S}} \eta \cdot \psi (\nabla_{\mathbf{x}}) \quad (2)$$

84 where $\nabla_{\mathbf{x}}$ is short for $\nabla_{\mathbf{x}} l(f_{\theta}(\mathbf{x}), y)$, η is the step size, ψ is a transformation function (e.g., *sgn* in
85 FGSM and *Identity* in FGM) and Π is the projection. Since such linearization attacks tend to be
86 trapped in the non-smooth vicinity of the data point, R+FGSM initializes a small *random* perturbation
87 to escape the vicinity and then applies FGSM. As a typical multi-step attacker, PGD^K [Madry et al.,
88 2018] can find better perturbations by K step gradient ascent:

$$\mathbf{x}^{(k)} = \Pi_{\mathbf{x}+\mathbb{S}} \left(\mathbf{x}^{(k-1)} + \eta \cdot \psi (\nabla_{\mathbf{x}^{(k-1)}}) \right) \quad (3)$$

89 3 Methodology

90 3.1 Motivation

91 The key of AT is to generate perturbations in the *inner maximization*. Strong perturbation helps to
92 improve robustness. It is generally accepted that the step K used to solve the *inner maximization*
93 correlates with the attacker’s ability to generate stronger perturbation. However, larger K leads to a
94 linear increase in training overhead. Wong et al. [2019] suggests that with appropriate step size tuning
95 and early stopping, one-step attackers yield models with the robustness that is comparable to much
96 more expensive multi-step attackers. It indicates that hyperparameters, such as random initialization,
97 step size, momentum, and early stopping, affect perturbation generation. From the perspective of
98 effectiveness and efficiency, it is valuable to further improve robustness by tuning the attacker to
99 strengthen perturbations.

100 However, manual tuning of perturbation generation for each example on-the-fly during training is
101 impractical. To address this problem, we propose an *efficient automated attacker* to boost adversarial
102 training by generating optimal perturbations on-the-fly during training.

103 3.2 Problem Formulation

104 Inspired by AutoML, we first design a general attacker space \mathcal{A} by referring to existing attackers.
105 Then, we employ an *automated attacker* parameterized by α to search in \mathcal{A} and further construct
106 an attacker against the example and the defense model $(\mathbf{x}, y, f_{\theta})$. We abbreviate the perturbation
107 generated by the constructed attacker as δ_{α} . Therefore, the goal of A² is to train a robust model using
108 the perturbation generated by the constructed attackers through a bilevel optimization problem:

$$\begin{aligned} & \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \in D} [l(f_{\theta}(\mathbf{x} + \delta_{\alpha^*}), y)] \\ & \text{s.t. } \alpha^* = \arg \max_{\alpha} \mathbb{E}_{(\mathbf{x}, y) \in D} [l(f_{\theta}(\mathbf{x} + \delta_{\alpha}), y)] \end{aligned} \quad (4)$$

109 On the attack side, we train α by SGD to make the defense model misclassify. On the defense
110 side, we use α^* to construct the best attacker for each example and then generate perturbations to
111 adversarially train the defense model.

112 3.3 Attacker Space

113 Revisiting most attackers, we find that the attacker can be viewed as a stack of one-step attackers
114 consisting of an attack method and a step size. Thus, as shown in Figure 1(a), we design a general
115 attacker space \mathcal{A} consisting of K -step cells. The k -th cell is denoted as $C^{(k)}$, which is a one-step
116 attacker consisting of the following two blocks:

117 **Perturbation Block O_p .** Typical attack methods (i.e., *FGM* and *FGSM*), attack methods with
118 momentum (i.e., *FGMM* and *FGSMM*), random perturbations (i.e., *Gaussian* and *Uniform*), and
119 the special *Identity* which enables the attacker to automatically early stop at a certain step like FAT
120 [Zhang et al., 2020];

121 **Step Size Block O_s .** $\{10^{-4} \cdot \eta, 10^{-3} \cdot \eta, 10^{-2} \cdot \eta, 10^{-1} \cdot \eta, \eta\}$, where η is a hyperparameter related to
122 the space of allowed perturbations \mathbb{S} .

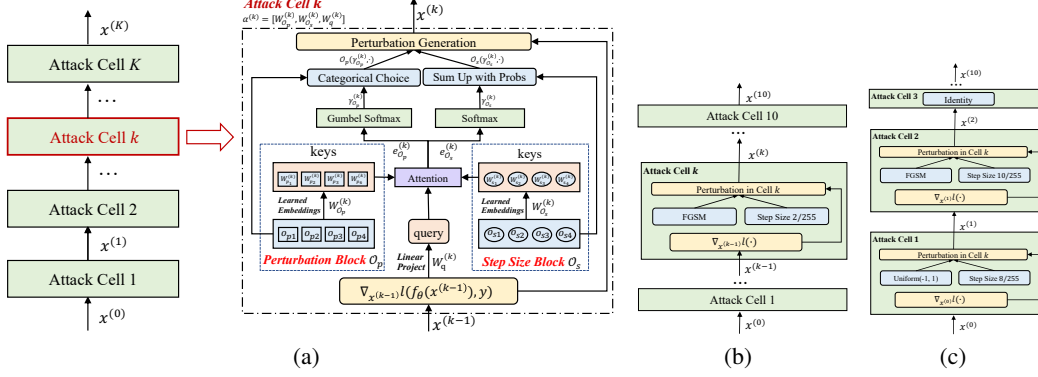


Figure 1: (a) Attacker Space of A^2 ; (b) PGD¹⁰: *FGSM* and a fixed step size 2/255 in each cell; (c) R+FGSM: *Gaussian* in the first cell, *FGSM* in the second cell and *Identity* in the other cells.

123 Each block O contains multiple operations. Let $o(\cdot)$ denote the operation, $\gamma^{(k)} = [\gamma_{O_p}^{(k)}, \gamma_{O_s}^{(k)}]$ denote
 124 the choice of operation in the k -th cell. The attack methods within $O_p^{(k)}$ are mutually exclusive. Thus,
 125 $\gamma_{O_p}^{(k)}$ is a one-hot vector. In contrast, the operations within $O_s^{(k)}$ are continuous. $\gamma_{O_s}^{(k)}$ is a normalized
 126 continuous vector, where each element represents the selection probability. To unify the categorical
 127 choice of attack methods and the probabilities over step sizes, the output of $O^{(k)}$ is expressed as a
 128 mixture based on $\gamma_{O}^{(k)}$:

$$\bar{O}(\gamma_{O}^{(k)}, \nabla_{\mathbf{x}^{(k-1)}}) = \sum_{o \in O^{(k)}} \gamma_o \cdot o(\nabla_{\mathbf{x}^{(k-1)}}) \quad (5)$$

129 where γ_o denotes the weight of the operation o in $\gamma_{O}^{(k)}$. Correspondingly, the one-step attacker of the
 130 k -th cell can be expressed as the joint of two blocks:

$$\bar{C}(\gamma^{(k)}, \nabla_{\mathbf{x}^{(k-1)}}) = \bar{O}_s(\gamma_{O_s}^{(k)}, \nabla_{\mathbf{x}^{(k-1)}}) \cdot \bar{O}_p(\gamma_{O_p}^{(k)}, \nabla_{\mathbf{x}^{(k-1)}}) \quad (6)$$

131 Moreover, the constructed attacker is a composition of attackers from each cell. In this way, we can
 132 cover common attackers in our space. For example, as shown in Figure 1(b), PGD^K is obtained by
 133 selecting *FGSM* in each perturbation block. R+FGSM in Figure 1(c) is a case of selecting *Gaussian*
 134 in the first cell, *FGSM* in the second cell and *Identity* in the other cells.

135 **Analysis of \mathcal{A}** Considering there exist 7 attack methods in O_p of each step, there are 7^K combinations
 136 of attack methods in the K -step attacker space. The exponential increasing combinations prevent the
 137 brute-force search. Moreover, the continuous step size O_s is also part of the attacker space. Thus, we
 138 propose A^2 to search for the best attacker in \mathcal{A} and generate adversarial perturbations efficiently.

139 3.4 Automated Attacker A^2

140 A^2 is used to construct the best attacker against $(\mathbf{x}, y, f_\theta)$ and its trainable parameters α include $W_{O_p}^{(k)}$,
 141 $W_{O_s}^{(k)}$, and $W_q^{(k)}$ where $k \in \{1, \dots, K\}$. As shown in Figure 1(a), we treat the current model and
 142 example as a query and the candidate operations as keys. Thus, the attention mechanism can be used
 143 to calculate the scores of operations within each block and the operations are selected based on their
 144 scores. Specifically, in the k -th cell, we take the gradient of the last step $\nabla_{\mathbf{x}^{(k-1)}} l(f_\theta(\mathbf{x}^{(k-1)}), y)$ as
 145 input and project it to a vector space as the query using $W_q^{(k)}$. Then, we use the trainable embedding
 146 table $W_o^{(k)}$ to convert the individual operations within $O^{(k)}$ to continuous keys. With the Scaled
 147 Dot-Product Attention [Vaswani et al., 2017], we compute the dot products of the query with each
 148 key as the score of the operation $o \in O$ in the k -th cell:

$$e_o^{(k)} = (\nabla_{\mathbf{x}^{(k-1)}} W_q^{(k)})^T W_o^{(k)} \quad (7)$$

149 **Perturbation Block.** The operations within $O_p^{(k)}$ are mutually exclusive. We sample an operation
 150 with the normalized scores as probabilities, i.e., $\gamma_{O_p}^{(k)} \sim \text{softmax}(e_{O_p}^{(k)})$.

151 **Step Size Block.** As the operations within $\mathcal{O}_s^{(k)}$ are continuous values, we sum up the individual step
 152 sizes with the normalized scores as weights. For $o_s \in \mathcal{O}_s^{(k)}$, the weight can be expressed as:

$$\gamma_{o_s}^{(k)} = \frac{\exp(e_{o_s}^{(k)})}{\sum_{o'_s \in \mathcal{O}_s^{(k)}} \exp(e_{o'_s}^{(k)})} \quad (8)$$

153 3.5 Training of Automated Attacker

154 As mentioned above, we train A^2 to minimize the following objective by gradient descent:

$$\alpha^* = \arg \min_{\alpha} -\mathbb{E}_{(\mathbf{x}, y) \in D} [l(f_{\theta}(\mathbf{x} + \delta_{\alpha}), y)] \quad (9)$$

155 However, as a result of constructing the attacker by sampling in each perturbation block, the gradient
 156 of the loss w.r.t $\gamma_{\mathcal{O}_p}$ is zero. To train α , we use the reparameterization trick [Kingma and Welling,
 157 2013] to transfer the randomness of sampling to the auxiliary noise and reformulate the objective
 158 function. For brevity, we omit the step index k .

159 Let $\gamma_{\mathcal{O}_p} = \phi(\kappa, e_{\mathcal{O}_p})$ be a differentiable transformation where κ is an auxiliary noise variable with
 160 independent marginal $p(\kappa)$. In A^2 , we sample noise from Gumbel Distribution, i.e., $\kappa \sim \text{Gumbel}(0)$
 161 [Gumbel, 1954], and use Gumbel Softmax [Jang et al., 2017] as ϕ to smoothly approximate the
 162 expectation of loss [Maddison et al., 2014]. Specifically, $\phi(\kappa, e_{\mathcal{O}_p}) = \text{softmax}((e_{\mathcal{O}_p} + \kappa)/\tau)$ where
 163 τ is the temperature parameter. When $\tau \rightarrow 0$, the generated samples have the same distribution as
 164 *one_hot*($\arg \max_{o_p \in \mathcal{O}_p} (e_{o_p} + \kappa_{o_p})$).

165 Using the reparameterization trick, we can now form MC (Monte Carlo) estimates of the expectation
 166 of A^2 's loss l for each example, which is differentiable, as follows:

$$\begin{aligned} & l(f_{\theta}(\mathbf{x} + \delta_{\alpha}), y) \\ &= \mathbb{E}_{\gamma_{\mathcal{O}_p} \sim \text{softmax}(e_{\mathcal{O}_p})} \left[l\left(f_{\theta}(\mathbf{x} + \bar{C}([\gamma_{\mathcal{O}_p}, \gamma_{\mathcal{O}_s}], \nabla_{\mathbf{x}}), y)\right) \right] \\ &= \mathbb{E}_{p(\kappa)} \left[l\left(f_{\theta}(\mathbf{x} + \bar{C}([\phi(\kappa, e_{\mathcal{O}_p}), \gamma_{\mathcal{O}_s}], \nabla_{\mathbf{x}}), y)\right) \right] \\ &\approx \frac{1}{M} \sum_{m=1}^M l\left(f_{\theta}\left(\mathbf{x} + \bar{C}([\phi(\kappa^{(m)}, e_{\mathcal{O}_p}), \gamma_{\mathcal{O}_s}], \nabla_{\mathbf{x}})\right), y\right) \end{aligned} \quad (10)$$

167 where M is the number of samples. In practice, $M = 1$ can achieve good performance. In this way,
 168 we reformulate the MC approximation in Equation (10) of l as the objective function \hat{l} .

169 Moreover, training α to convergence in each epoch can be prohibitive due to the expensive inner
 170 maximization in Equation (4). We use a simple approximation scheme following the common
 171 methods [Finn et al., 2017, Liu et al., 2018]:

$$\alpha^* \approx \alpha + \xi \nabla_{\alpha} \mathbb{E}_{(\mathbf{x}, y) \in D} [\hat{l}(f_{\theta}(\mathbf{x} + \delta_{\alpha}), y)] \quad (11)$$

172 where α denotes the current weights of the attacker and ξ is the learning rate.

173 3.6 Framework of Adversarial Training with A^2

174 The overall procedure is shown in Algorithm 1. As in normal adversarial training, we generate
 175 perturbations in K steps every batch and update the model parameters. The key difference is in Line
 176 7. Benefiting from a parameterized automated attacker, we tune the discrete attack methods and
 177 continuous step sizes to generate adversarial perturbations. After optimizing the model parameters, we
 178 use Equation (11) to update α as an approximation to α^* . Since A^2 focus on the *inner maximization*,
 179 it can be compatible with most adversarial training methods. For example, it is flexible to use the loss
 180 function of TRADES or MART for *outer minimization* in Line 10 (i.e., TRADES- A^2 and MART- A^2),
 181 or include early stopping in Line 5~9 as FAT.

182 4 Experiments

183 We conduct extensive experiments on public datasets to answer the following questions: 1) Can A^2
 184 generate stronger adversarial perturbations? 2) How effective is the adversarial training with A^2 ?

Algorithm 1 Adversarial Training with Automated Attacker (AT-A²)

Input: Training examples D , perturbation bound ϵ , the number of attack steps K

```

1: Initialize  $\theta, \alpha$ ;
2: for epoch = 1, ...,  $N_{ep}$  do
3:   for minibatch  $(\mathbf{X}, Y) \subset D$  do
4:      $\mathbf{X}^{(0)} \leftarrow \mathbf{X}$ ;
5:     for  $k = 1, \dots, K$  do
6:       Calculate the gradient  $\nabla_{\mathbf{X}^{(k-1)}}$ ;
7:       Construct  $\delta_{\alpha}^{(k)} \in \mathbb{S}$  according to  $\nabla_{\mathbf{X}^{(k-1)}}$  by  $g_{\alpha}$ ;
8:        $\mathbf{X}^{(k)} = \mathbf{X}^{(k-1)} + \delta_{\alpha}^{(k)}$ ;
9:     end for
10:    Update  $\theta$  with  $\nabla_{\theta} \sum_{(x,y)} l(f_{\theta}(\mathbf{x}^{(K)}), y)$ ;
11:    Update  $\alpha$  by Equation (11);
12:  end for
13: end for
  
```

185 3) Is A² robust to hyperparameters? All experiments are run using GeForce RTX 3090 (GPU) and
 186 Intel(R) Xeon(R) Silver 4210 (CPU) instances.

187 **4.1 Effectiveness of Automated Attacker (RQ1)**

188 In this part, we fix the model f_{θ} , train the automated attacker alone and investigate whether A² can
 189 generate more powerful perturbations compared to the commonly used PGD.

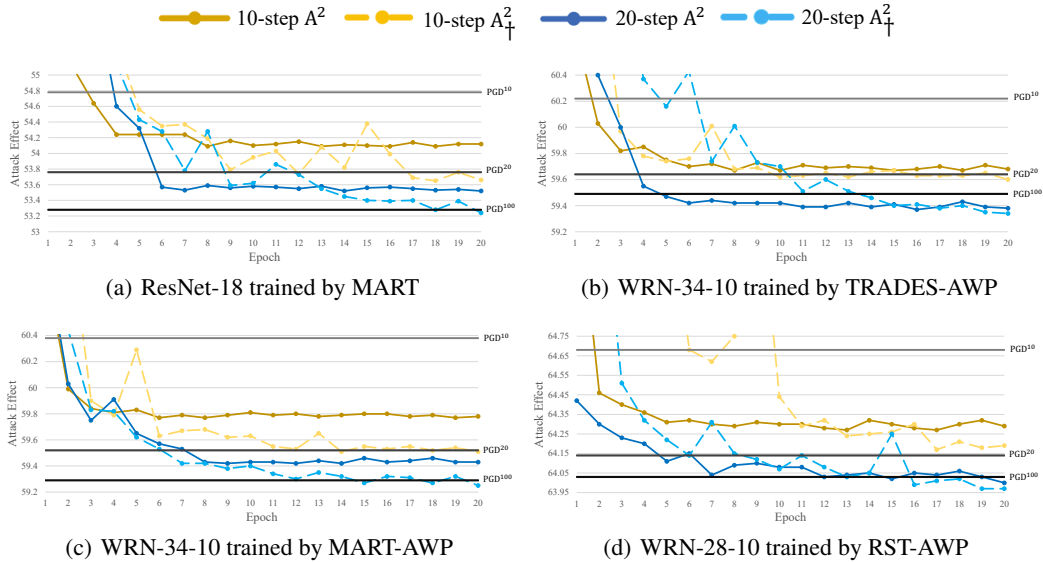


Figure 2: Effect of adversarial perturbations generated by A² with the training epoch.

190 **Experimental Settings.** To demonstrate the generality, we choose different neural architectures
 191 (i.e., ResNet-18, WRN-34-10, and WRN-28-10) trained on CIFAR-10 by various AT methods:
 192 TRADES, MART, RST, and AWP. All trained models are open-source checkpoints. We choose
 193 PGD^K with a random start $\delta^{(0)} \sim Uniform(-\epsilon, \epsilon)$ as baseline. All attacks are L_{∞} -bounded with a
 194 total perturbation scale of $\epsilon = 8/255$. Since more attack steps generally improve the attack effect,
 195 we compare the automated attacker with PGD at different steps. Moreover, we try different step
 196 size blocks, i.e., η , in A²: 1) $\eta = 2/255$, which is the same as the setting of step size in PGD; 2)
 197 $\eta = 8/255$, which is indicated by A²_‡ and allows A² to search the whole ϵ bound each step. A² is

Table 1: Comparison of attack effects (% , the lower the better) of multi-step PGD and A^2 in robust models. We run each method 5 times and show the average. The standard deviations are omitted as they are very small. The architecture of all defense models is WideResNet, except for MART whose architecture is ResNet-18.

Defense	Natural	10-step			20-step			PGD ¹⁰⁰
		PGD	A^2	A^2_{\dagger}	PGD	A^2	A^2_{\dagger}	
MART ⁰	83.07	54.78	54.09	53.65	53.76	53.52	53.24	53.28
TRADES-AWP ¹	85.36	60.22	59.67	59.60	59.64	59.38	59.34	59.49
MART-AWP ¹	85.60	60.38	59.76	59.51	59.52	59.42	59.25	59.29
RST-AWP ¹	88.25	64.68	64.27	64.17	64.14	64.02	63.97	64.03

198 trained using Adam [Kingma and Ba, 2015] with learning rate 10^{-3} , weight decay 10^{-2} and other
 199 default hyperparameters.

200 **Attack Effect.** The training process of A^2 is shown in Figure 2. We take the first 20 epochs of the
 201 attack effects, compare them with PGD, and observe whether A^2 converges. In the early training
 202 stage, the random combinations of attack operations are much less effective. After 10~20 epochs,
 203 the effect of generated attacks is much more stable and effective. In practical automated adversarial
 204 training, the model and the automated attacker are trained iteratively. The fast convergence of A^2
 205 ensures that the generated perturbations are strong enough. In addition, a larger η achieves better
 206 attacks. As the steps increase, the effect diminishes and the training may fluctuate.

207 Table 1 reports the attack effects of PGD and A^2 with different steps. In the comparison of A^2 and
 208 PGD with different steps $K \in \{10, 20\}$, A^2 stably outperforms PGD. A^2_{\dagger} constructs stronger attacks
 209 in the expanded search space. With the increase of steps, A^2 is more effective due to the combination
 210 of attack methods and the automated step size tuning. Due to the diminishing marginal effect, PGD¹⁰⁰
 211 v.s. PGD²⁰ achieves less improvement than PGD²⁰ v.s. PGD¹⁰. At 1/5 of the cost, the 20-step A^2_{\dagger}
 212 finds better attacks using the optimized α than PGD¹⁰⁰. In summary, Table 1 verifies that A^2 stably
 213 outperforms PGD for the same step, and obtains better attacks compared to PGD, whose step size
 214 and attack method are fixed, with significantly lower cost.

215 **Overhead Analysis.** The overhead of A^2 is not significant compared to PGD. Both methods are
 216 close in terms of clock time. For WRN-34, PGD takes 19.75/147.09/287.76 seconds to generate
 217 1/10/20 step attacks respectively. It demonstrates that more inner steps lead to a linear increase in
 218 time. Meanwhile, A^2 takes 157.61/302.51 seconds to generate the 10/20 step attack respectively. The
 219 main overhead remains in the forward computation and backward propagation of the defense model.
 220 Moreover, Section A.3 in Appendix shows that the total parameter size of A^2 is also acceptable.

221 4.2 Effectiveness of Adversarial Training with A^2 (RQ2)

222 In this part, we evaluate the robustness of our proposed AT- A^2 on different datasets against white-box
 223 and ensemble attacks. To verify that the stronger attacks generated by A^2 on-the-fly during training
 224 can improve robustness, we consider various adversarial training methods (i.e., AT, TRADES, MART,
 225 and AWP) without additional data across different datasets.

226 **Benchmark.** We conduct experiments on the baseline AT and the SOTA AWP with A^2 across three
 227 benchmark datasets to verify the generalization of A^2 . We follow the settings in AWP: PreActResNet-
 228 18 trained for 200 epochs, $\epsilon = 8/255$ and $\gamma = 10^{-2}$ for AWP. The step size is 1/255 for SVHN and
 229 2/255 for CIFAR-10 and CIFAR-100. For AT and AWP, the attacker used in training is PGD¹⁰.
 230 The 10-step A^2 is trained with the same setting as in RQ1. PGD²⁰ is used for testing, and the test
 231 robustness is reported in Table 2. It shows that A^2 , as a component focusing on the *inner maximization*,
 232 achieves better results on most datasets. Moreover, A^2 is generic and can boost the robustness of both
 233 baseline and SOTA AT methods.

⁰<https://github.com/YisenWang/MART>

¹<https://github.com/csdongxian/AWP>

²<https://github.com/zjfhart/Friendly-Adversarial-Training>

Table 2: Test robustness (% , the higher the better) using PreActResNet-18 under L_∞ threat model ("Best" means the highest robustness while "Last" means the robustness at the last epoch). Std. of 5 runs is omitted due to being small.

Defense	SVHN		CIFAR-10		CIFAR-100	
	Best	Last	Best	Last	Best	Last
AT	53.36	44.49	52.79	44.44	27.22	20.82
AT-A ²	56.76	44.75	52.96	44.59	28.14	20.28
AWP	59.12	55.87	55.39	54.73	30.71	30.28
AWP-A ²	61.42	58.45	55.71	55.31	31.36	30.73

Table 3: Test robustness (% , the higher the better) on CIFAR-10 using WRN-34-10 under L_∞ threat model ("Natural" denotes the accuracy on nature examples, and other columns indicate the accuracy on adversarial examples generated by different attacks). Std. of 5 runs is omitted due to being small.

Defense	Natural	FGSM	PGD ²⁰	CW _∞	AutoAttack
AT	87.30	56.10	52.68	50.73	47.04
AT-A ²	84.54	63.72	54.68	51.17	48.36
TRADES	84.65	61.32	56.33	54.20	53.08
TRADES-A ²	85.54	65.93	59.84	56.61	55.03
MART	84.17	61.61	57.88	54.58	51.10
MART-A ²	84.53	63.73	59.57	54.66	52.38
AWP	85.57	62.90	58.14	55.96	54.04
AWP-A ²	87.54	64.70	59.50	57.42	54.86

234 **Robustness on WideResNet.** Furthermore, we train WRN-34-10 on CIFAR-10 with various AT
 235 methods (i.e., AT, TRADES, MART, and AWP) following their original papers and open-source
 236 codes². All defense models are trained using SGD with momentum 0.9, weight decay 5×10^{-4} , and
 237 an initial learning rate of 0.1 that is divided by 10 at the 50%-th and 75%-th epoch. Except for 200
 238 epochs in AWP, other AT methods train the model for 120 epochs. Simple data augmentations (i.e.,
 239 32x32 random crop with 4-pixel padding and random horizontal flip) are applied.

240 For white-box attack, we test FGSM, PGD²⁰ and CW_∞ [Carlini and Wagner, 2017]. In addition,
 241 we test the robustness against the standard AutoAttack [Croce and Hein, 2020], which is a strong
 242 and reliable attacker to verify the robustness via an ensemble of diverse parameter-free attacks
 243 including three white-box attackers and a black-box attacker. Table 3 shows that A² reliably boosts
 244 AT variants against white-box and ensemble attacks. This verifies that A² is general for AT and
 245 improves adversarial robustness reliably rather than gradient obfuscation or masking.

246 Additionally, given the nature examples, AT performs better than AT-A². The main reason is that
 247 A² generates stronger perturbation for better robustness, which decreases the accuracy (i.e., 84.54).
 248 Many works (e.g., TRADES, MART, and AWP) use regularization to achieve the trade-off between
 249 robustness and accuracy. The regularization is also used to optimize the automated attacker. Thus, for
 250 other AT methods in Table 3, combining A² can achieve higher accuracy. Moreover, for WRN-34,
 251 the training time of AWP-A² is 970 s/epoch while the training time of AWP is 920 s/epoch. Thus, the
 252 additional overhead of A² is not significant.

253 4.3 Hyperparameters of A² (RQ3)

254 The hyperparameters of A² include the training hyperparameters and the design of the attacker space.
 255 The comparison of the attack effect with different hyperparameters is shown in Table 4. Overall, A²
 256 is robust to hyperparameters and performs better than PGD¹⁰ and closely to PGD²⁰.

257 **Training of A².** The effect of attacks with different learning rates ξ is shown in the middle two rows
 258 of Table 4. Although Adam uses a dynamic learning rate, an excessive initial learning rate (i.e., 10^{-2})
 259 leads to sub-optimal.

Table 4: Comparison of attack effects (% , the lower the better) of 10-step A^2 with different hyperparameters (The 10-step is omitted, $A^2_{p,q}$ is short for training the attacker using learning rate $\xi = 10^{-p}$ with the step size block $\eta = q/255$,). Std. of 5 runs is omitted due to being small.

Attack	TRADES-AWP ¹	MART-AWP ¹	RST-AWP ¹
PGD ¹⁰	60.22	60.38	64.68
PGD ²⁰	59.64	59.52	64.14
$A^2_{3,2}$	59.67	59.76	64.27
$A^2_{2,2}$	59.93	59.87	64.34
$A^2_{4,2}$	59.76	59.78	64.29
$A^2_{3,5}$	59.49	59.62	64.11
$A^2_{3,8}$	59.53	59.53	64.17

260 **Attacker Space.** The influence of the attack step K has been investigated in RQ1. As shown in the
 261 last two rows of Table 4, a larger step size η increases the effectiveness of A^2 . However, as shown in
 262 the training curve in Figure 2, larger η introduces instability in the training of the attacker.

263 5 Related Work

264 5.1 Adversarial Learning

265 Many recent works [Goodfellow et al., 2015, Carlini and Wagner, 2017, Croce and Hein, 2020] have
 266 shown that DNNs are vulnerable to adversarial examples. Various defense strategies and models have
 267 been proposed to deal with the threat of adversarial examples. However, as proved in C&W [Carlini
 268 and Wagner, 2017], many works mistake gradient obfuscation or masking for adversarial robustness.
 269 AT [Madry et al., 2018] formulates a class of adversarial training methods for solving a saddle point
 270 problem (i.e., Equation (1)) and improves robustness reliably.

271 Based on AT, many works [Zhang et al., 2019, Wang et al., 2019, Wu et al., 2020] focusing on
 272 the components of *outer minimization* are introduced to further enhance performance. The *inner*
 273 *maximization* is also the goal of the adversarial attack, where l is the 0-1 loss. Many works, e.g.,
 274 FGSM [Goodfellow et al., 2015], C&W [Carlini and Wagner, 2017] and AutoAttack [Croce and Hein,
 275 2020], have been proposed to attack DNNs and facilitate the development of adversarial training.

276 5.2 Automated Machine Learning

277 AutoML [Bergstra et al., 2011, Zoph and Le, 2017, Liu et al., 2018, Cubuk et al., 2019] aims to
 278 automate the parts of the machine learning pipeline that require expert solutions. For a particular
 279 domain, it is common practice to summarize a large search space of parameters and configurations
 280 based on expert experience and search for the optimal solutions using methods such as black-box
 281 optimization. During the search process, a certain metric is required to evaluate each solution. The
 282 same idea can be applied to adversarial learning. A^3 [Yao et al., 2021], which is also closely related
 283 to AutoML, automatically discovers an effective attacker on a given model.

284 6 Conclusion

285 In this work, we proposed A^2 , to the best of our knowledge, the first adversarial training method which
 286 focuses on automated perturbation generation. In A^2 , the attacker space is designed by summarizing
 287 the existing perturbations. Moreover, the parameterized automated attacker leverages the attention
 288 mechanism to choose the discrete attack method and the continuous step size and further generates
 289 adversarial perturbations. During training, the one-step approximation of the optimal automated
 290 attacker is used to generate the optimal perturbations on-the-fly for the model. The experimental
 291 results show that A^2 generates stronger attacks with low extra cost and boosts the robustness of
 292 various AT methods reliably.

293 For future work, we plan to add the target loss of the *inner maximization* to the attacker space. We
 294 also plan to apply A^2 to enhance adversarial training for Natural Language Processing.

295 **References**

- 296 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
297 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
298 2016.
- 299 Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep
300 bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, 2019.
- 301 Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorization-
302 machine based neural network for ctr prediction. In *Proc. of IJCAI*, 2017.
- 303 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *BMVC*, 2016.
- 304 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
305 Theoretically principled trade-off between robustness and accuracy. In *Proc. of ICML*, 2019.
- 306 Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving
307 adversarial robustness requires revisiting misclassified examples. In *Proc. of ICLR*, 2019.
- 308 Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust general-
309 ization. *Advances in Neural Information Processing Systems*, 2020.
- 310 Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, Percy Liang, and John C Duchi. Unlabeled data
311 improves adversarial robustness. In *Proc. of ICONIP*, 2019.
- 312 Sven Gowal, Chongli Qin, Jonathan Uesato, Timothy Mann, and Pushmeet Kohli. Uncovering
313 the limits of adversarial training against norm-bounded adversarial examples. *arXiv preprint*
314 *arXiv:2010.03593*, 2020.
- 315 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
316 Towards deep learning models resistant to adversarial attacks. In *Proc. of ICLR*, 2018.
- 317 Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
318 examples. In *Proc. of ICLR*, 2015.
- 319 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
320 of diverse parameter-free attacks. In *International conference on machine learning*, 2020.
- 321 Chengyuan Yao, Pavol Bielik, Petar Tsankov, and Martin Vechev. Automated discovery of adaptive
322 attacks on adversarial defenses. *arXiv preprint arXiv:2102.11860*, 2021.
- 323 Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. In
324 *Proc. of ICLR*, 2019.
- 325 Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *Proc. of*
326 *ICLR*, 2017.
- 327 Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *Proc.*
328 *of ICLR*, 2018.
- 329 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. 2017.
- 330 Takeru Miyato, Andrew M Dai, and Ian Goodfellow. Adversarial training methods for semi-supervised
331 text classification. In *Proc. of ICLR*, 2017.
- 332 Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankan-
333 halli. Attacks which do not kill training make adversarial learning stronger. In *ICML*, 2020.
- 334 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
335 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
336 *processing systems*, 2017.
- 337 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
338 *arXiv:1312.6114*, 2013.

- 339 Emil Julius Gumbel. *Statistical theory of extreme values and some practical applications: a series of*
340 *lectures*. US Government Printing Office, 1954.
- 341 Chris J Maddison, Daniel Tarlow, and Tom Minka. A* sampling. *Advances in Neural Information*
342 *Processing Systems*, 2014.
- 343 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of
344 deep networks. In *Proc. of ICML*, 2017.
- 345 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. of ICLR*,
346 2015.
- 347 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017*
348 *IEEE Symposium on Security and Privacy (SP)*, 2017.
- 349 James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter
350 optimization. *Advances in neural information processing systems*, 2011.
- 351 Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment:
352 Learning augmentation strategies from data. In *Proc. of CVPR*, 2019.
- 353 Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting
354 adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision*
355 *and pattern recognition*, pages 9185–9193, 2018.
- 356 Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L Yuille.
357 Improving transferability of adversarial examples with input diversity. In *Proceedings of the*
358 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2730–2739, 2019.
- 359 Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Evading defenses to transferable adversarial ex-
360 amples by translation-invariant attacks. In *Proceedings of the IEEE/CVF Conference on Computer*
361 *Vision and Pattern Recognition*, pages 4312–4321, 2019.

362 Checklist

- 363 1. For all authors...
- 364 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
365 contributions and scope? [Yes]
- 366 (b) Did you describe the limitations of your work? [Yes]
- 367 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 368 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
369 them? [Yes]
- 370 2. If you are including theoretical results...
- 371 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 372 (b) Did you include complete proofs of all theoretical results? [N/A]
- 373 3. If you ran experiments...
- 374 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
375 mental results (either in the supplemental material or as a URL)? [Yes]
- 376 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
377 were chosen)? [Yes]
- 378 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
379 iments multiple times)? [Yes] The standard deviations are omitted as they are very
380 small
- 381 (d) Did you include the total amount of compute and the type of resources used (e.g., type
382 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix.
- 383 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 384 (a) If your work uses existing assets, did you cite the creators? [Yes]

- 385 (b) Did you mention the license of the assets? [Yes]
 386 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 387
 388 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 389 using/curating? [N/A]
 390 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 391 information or offensive content? [N/A]
 392 5. If you used crowdsourcing or conducted research with human subjects...
- 393 (a) Did you include the full text of instructions given to participants and screenshots, if
 394 applicable? [N/A]
 395 (b) Did you describe any potential participant risks, with links to Institutional Review
 396 Board (IRB) approvals, if applicable? [N/A]
 397 (c) Did you include the estimated hourly wage paid to participants and the total amount
 398 spent on participant compensation? [N/A]

399 A Details in A²

400 A.1 Unify Magnitude of Perturbations

401 Perturbations generated by different operations in \mathcal{O}_p have different magnitudes and thus require
 402 different magnitudes of step size for different o_p . For example, *FGSM* generates perturbations
 403 with elements belonging to $\{-1, 0, 1\}$, while the perturbation generated by *FGM* is usually in the
 404 magnitude of 10^{-3} . Obviously, they cannot use the same step size. To have a uniform effect of the
 405 step size block, we normalize the magnitude of other generated perturbations to be the same as *FGSM*
 406 (i.e., $\delta_{o_p} = \delta_{o_p} \cdot \frac{\|\delta_{FGSM}\|}{\|\delta_{o_p}\|}$). In this way, we find that other attack methods such as *FGM* can achieve
 407 good results with the same step size as *FGSM*.

408 A.2 Temperature Parameter in *Softmax*

409 Since there is an order of magnitude difference in step size operations, the larger step size with the
 410 same score will dominate the output. For example, $0.7 \cdot 10^{-2}\eta + 0.3 \cdot \eta \approx 0.3\eta$. The output of the
 411 step size block is dominated by the operation η , despite the greater weight of $10^{-2} \cdot \eta$. To alleviate
 412 the problem, we use the temperature parameter τ in *softmax* to sharpen the distribution:

$$\gamma_{o_s}^{(k)} = \frac{\exp(e_{o_s}^{(k)}/\tau)}{\sum_{o' \in \mathcal{O}_s} \exp(e_{o'}/\tau)} \quad (12)$$

413 where o_s is an operation in \mathcal{O}_s , and e_{o_s} is its attention score. Through experiments, we set $\tau = 0.1$ to
 414 distinguish the preference for the step size in most cases.

415 A.3 Overhead of A²

416 Let the number of steps be K , the number of operations be $|\mathcal{O}|$, the image size be $W \times H$ and
 417 the embedding size be E . The number of the attacker’s parameter is $\mathbf{O}(K \cdot E \cdot (W \times H + |\mathcal{O}|))$.
 418 Specifically, the number of parameters for the attacker is 7873280, which is 17% of the model’s
 419 parameters (i.e., 46160474). In each batch, there is only 1 forward calculation of all cells with 1
 420 backpropagation. In comparison, the model requires K forward calculations with backpropagation.
 421 Therefore, the additional computational overhead from the attacker is not significant in terms of the
 422 number of parameters and computations.

423 Moreover, PGD and A² are close in terms of clock time. For WRN-34, PGD takes
 424 19.75/147.09/287.76 seconds to generate 1/10/20 step attacks respectively. It demonstrates that
 425 more inner steps lead to a linear increase in time. Meanwhile, A² takes 157.61/302.51 seconds to
 426 generate the 10/20 step attack respectively. The main overhead remains in the forward computation
 427 and backward propagation of the defense model. For WRN-34, the training time of AWP-A² is 970
 428 s/epoch while the training time of AWP is 920 s/epoch.

429 In summary, the additional overhead of A² is not significant.

430 **A.4 Why No Mixture in O_p**

431 Like most NAS methods in AutoML, the discrete selection in the perturbation block is more inter-
 432 pretable and robust (e.g., L1-Norm for feature selection and single path in NAS) than the mixture
 433 over possible solutions. Moreover, the mixture will incur more computational overhead and 7 times
 434 memory overhead due to 7 operations in O_p . Figure 3 shows an example of the generated attack on
 435 CIFAR-10, which can be migratable.

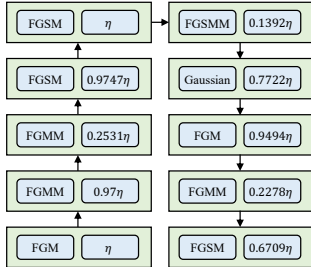


Figure 3: Example of generated attack on CIFAR-10.

436 **B Addition Experiments**

437 **B.1 Why use FGSM-based PGD in RQ1.**

438 There are multiple single-step attack methods in O_p for stacking as PGD, e.g., FGM-based PGD
 439 and FGSM-based PGD. The experimental results of the attack effect of PGD based on these attack
 440 methods demonstrate that FGSM-based PGD outperforms the stacking of other operations. Thus, we
 441 choose FGSM-based PGD with a random start $\delta^{(0)} \sim Uniform(-\epsilon, \epsilon)$ as a baseline for comparison
 442 with the automated attacker.

443 **B.2 Number of samples M in MC Approximation.**

444 M is an important hyperparameter that dictates the quality of MC approximation and the training
 445 overhead. We test the cases with $M \in \{1, 2, 5\}$ and achieve similar performance. Thus, we set M to 1
 446 and achieve good results with a significantly lower overhead.

447 **B.3 Generality of A^2 in White-Box Attacks**

Table 5: Comparison of attack effects on CIFAR-10 (% , the lower the better) of PGD-based and CW_∞ -based attacks. The architecture of all defense models is WideResNet, except for MART whose architecture is ResNet-18.

	MART	TRADES-AWP	MART-AWP	RST-AWP
Natural	83.07	85.36	85.60	88.25
PGD ²⁰	53.76	59.64	59.52	64.14
PGD ²⁰ -A ²	53.24	59.34	59.25	63.97
CW_∞	49.97	57.07	56.44	61.82
CW_∞ -A ²	49.82	56.98	55.81	61.30

448 In this part, we investigate whether A^2 is general to white-box attacks. As a more powerful attack
 449 method, CW_∞ -based attacks [Carlini and Wagner, 2017] stably outperform PGD-based attacks. For
 450 comparison with CW_∞ , we propose a variant of A^2 that uses CW_∞ loss to generate perturbations and
 451 denote it as CW_∞ -A². The results in Table 5 show that A^2 is general and can improve the attack effect
 452 of PGD and CW_∞ by combining attack methods and tuning the step size. Moreover, the additional
 453 overhead of A^2 is 5% to 10%, which is a rather acceptable trade-off.

454 **B.4 Robustness Against Transferable Black-Box Attacks**

455 We investigate the robustness of A^2 against transferable black-box attacks. Table 6 provides test
 456 robustness on CIFAR-10 using ResNet-18. We adopt three transferable black-box attack methods:
 457 MI (momentum = 1) [Dong et al., 2018], DI [Xie et al., 2019], and TI [Dong et al., 2019]. The
 458 transferable attacks are generated by an ensemble of the above methods on three surrogate pre-trained
 459 models ¹: IncV3 (InceptionV3), VGG19, and DN201 (DenseNet201). Table 6 shows that AT boosts
 460 the robustness against transferable black-box attacks, and A^2 can further improve the adversarial
 461 robustness.

Table 6: Test robustness (% , the higher the better) on CIFAR-10 using ResNet-18 against transferable black-box attacks.

	MI+DI+TI			
	IncV3	VGG19	DN201	PGD ²⁰
ResNet-18	16.12	7.37	5.35	0.02
ResNet-18-AT	61.98	60.81	59.63	52.79
ResNet-18-AT- A^2	62.79	61.85	60.28	52.96

462 **B.5 A Closer Look at Selected Attacks**

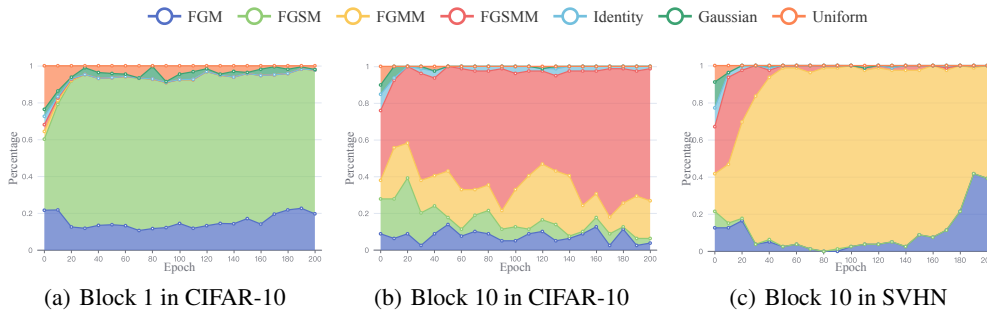


Figure 4: Distribution of attacks selected by perturbation blocks of A^2 .

463 We analyze the selected attacks from the perspective of perturbation blocks with different steps and
 464 datasets.

465 The first and final perturbation blocks of 10-step A^2 in CIFAR-10 are chosen for analysis. Figure 4
 466 shows the distribution of selected attacks of different perturbation blocks.

- 467 • **Perturbation Block 1:** A^2 tends to choose *FGM*, *FGSM*, and partially random methods as
 468 initialization in the first step. The momentum-based attack methods are quickly discarded
 469 as the gradient of the previous step is absent. *FGSM* is chosen more frequently due to its
 470 stronger attack on both foreground and background.
- 471 • **Perturbation Block 10:** The optimization of the victim model leads to changes in the
 472 distribution of selected attacks in the last block. In the early stage of training, the victim
 473 model is vulnerable. A^2 retains the diversity and plays the role of friendly attackers like
 474 FAT Zhang et al. [2020]. At the end of training, A^2 prefers the momentum-based attacks
 475 (i.e., *FGSM* and *FGMM*).

476 From the perspective of datasets, SVHN and CIFAR-10 prefer different attack methods. As shown in
 477 Figure 4(c), SVHN discards *FGSM*, which is most frequently used in CIFAR-10, and pays more
 478 attention to *FGMM*. Moreover, SVHN rarely uses *Identity* compared with CIFAR-10 as its higher
 479 robustness accuracy requires more powerful perturbations.

480 In summary, A^2 's preference for selecting attacks in blocks varies according to the block step, dataset,
 481 and victim model.

¹https://github.com/huyvnphan/PyTorch_CIFAR10