# Unifying Gradients to Improve Real-World Robustness for Deep Networks

YINGWEN WU, Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, China

SIZHE CHEN, Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, China

KUN FANG, Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, China

XIAOLIN HUANG, Institute of Image Processing and Pattern Recognition and the MOE Key Laboratory of System Control and Information Processing, Shanghai Jiao Tong University, China

The wide application of deep neural networks (DNNs) demands an increasing amount of attention to their real-world robustness, *i.e.*, whether a DNN resists black-box adversarial attacks, among which score-based query attacks (SQAs) are most threatening since they can effectively hurt a victim network with the only access to model outputs. Defending against SQAs requires a slight but artful variation of outputs due to the service purpose for users, who share the same output information with SQAs. In this paper, we propose a real-world defense by Unifying Gradients (UniG) of different data so that SQAs could only probe a much weaker attack direction that is similar for different samples. Since such universal attack perturbations have been validated as less aggressive than the input-specific perturbations, UniG protects real-world DNNs by indicating attackers a twisted and less informative attack direction. We implement UniG efficiently by a Hadamard product module which is plug-and-play. According to extensive experiments on 5 SQAs, 2 adaptive attacks and 7 defense baselines, UniG significantly improves real-world robustness without hurting clean accuracy on CIFAR10 and ImageNet. For instance, UniG maintains a model of 77.80% accuracy under 2500-query Square attack while the state-of-the-art adversarially-trained model only has 67.34% on CIFAR10. Simultaneously, UniG outperforms all compared baselines in terms of clean accuracy and achieves the smallest modification of the model output. The code is released at https://github.com/snowien/UniG-pytorch.

CCS Concepts: • **Computing methodologies** → **Computer vision**; **Neural networks**; • **Security and privacy** → Security services.

Additional Key Words and Phrases: Black-box Adversarial Attack, Practical Adversarial Defense

## 1 INTRODUCTION

Deep neural networks (DNNs) have been revealed to be vulnerable to adversarial examples (AEs), which can mislead models into incorrect predictions by imperceptible perturbations on inputs [21, 48, 60]. Such sensitivity poses a threat to real-world applications of DNNs, since an attacker only needs the same information as the user, namely the model output, to generate valid AEs. Considering that the inner information of the network, such as parameters or training datasets, is hidden from attackers in practical scenarios, the query-based attack [1, 3, 7, 11, 16, 23, 27, 28], which only requires the model output to hurt a victim model, deserves more attention in the field of real-world robustness compared to white-box attacks [9, 21] and transferable attacks [14, 15, 26]. Moreover, since plentiful applications output both prediction results and probabilities [6, 6, 10] to users for better judgement, score-based query attacks (SQAs) are more threatening because of their effectiveness and feasibility, compared to decision-based query attacks (DQAs), which needs unreasonable number of queries to attack.

Since SQAs are based on output probabilities to generate AEs, current defenses are all designed to alter the probabilities, either directly or indirectly, to resist attackers. For instance, adversarial training (AT), the most popular defense, uses on-the-fly AEs as training data to obtain a robust model whose output probabilities are always under-confident [22, 33, 36, 49]. Randomness injection (RI) is another effective defense against SQAs, which inject randomness into inputs [17, 44, 54], parameters [20, 25, 35], or features [34] to ultimately change probabilities with random noises to confound attackers. In contrast to RI, denoising methods [2, 37, 40] pre-process inputs to mitigate adversarial noises and reconstruct natural images, such that the output probability of adversarial queries can be modified to be consistent with the score of the corresponding clean query. Additionally, dynamic defenses [51, 52], which optimize model parameters at inference time to adapt attacks, also try to keep the same output probabilities with clean data when adversarial queries come. Although above defenses could mitigate SQAs, however, they hurt clean accuracy, which is a common phenomenon called accuracy-robustness trade-off [43, 45, 50, 53, 59]. An intuitive explanation to this trade-off is that high clean accuracy requires the use of detailed features, while they are the inducement of vulnerability [50]. From distribution perspective, natural images and adversarial examples belong to different distributions and thus cannot be fitted well simultaneously, which is simply proved by AT. In addition to the degradation of clean accuracy, we discover that these defenses inevitably affects the output probabilities of clean images, which seriously influences downstream tasks to make reasonable decisions, $e.g.$, the detection network needs to report accurate confidence to the center controller to avoid erroneous decisions on the object with low confidence. Therefore, in this paper, we aim at keeping accuracy and probabilities of clean data and meanwhile changing output probabilities of adversarial queries to simultaneously serve users and resist attackers.

Achieving the above goal is difficult because that both users and attackers share the same information, $i.e.$, the model output probability, while we need to keep the probability of clean data and change the probability of adversarial queries in the condition of unknown input types. In spite of this, our chance lies in the fact that users only ask for outputs, while SQAs concentrate on the change in output indicated by different queries, which implies gradient information used to attack. Therefore, we propose a defense that explicitly changes the gradient information contained in the output of consecutive queries. Through slight but designed modifications on outputs, which guarantee the service to the user, the gradient information contained in the output is perturbed to be an elaborate direction which is less aggressive. We choose the direction of the universal attack perturbation (UAP, [58]) here, which is consistent for different inputs and hereby less threatening than normal adversarial noises that are image-specific. Previous studies have proved the gentleness of UAP empirically [5, 56–58]. As a result, even through the attacker mines the attack direction from these outputs, their attack trajectory will be tricked away from the vulnerable adversarial direction and induced into the weaker path we have designated.

To confuse SQAs towards the UAP direction, the modification on outputs needs to be designed carefully. A natural idea is that we use a gradient unification loss, which constrains the gradient of different inputs to be the same, to optimize the output changes. The proposed method is then called unifying gradients (UniG) approach. Considering computational overhead, we choose to unify the gradient of features instead of inputs to efficiently calculate the second derivative in practice, as shown in Figure 1(a). We insert a Hadamard product module $A$ into a pre-trained DNN, where we want $\hat{f} \approx f$ for user friendliness and $\hat{g} \neq g$ but rather a UAP direction for adversarial robustness. The specific calculation of $\hat{f}, \hat{g}$ is that $\hat{f} = A \circ f$, $\hat{g}_i = g_i \circ A_i$, where $A \in R^{b \times d}$, $A_i \in R^d$ ($b$ for batch size and $d$ for feature dimension) is the module parameter and $g_i \in R^d$ is the feature gradient of the $i$-th input. According to the design of $\hat{f}$ and $\hat{g}$, we constrain each element of $A$ to be close to one to ensure slight forward modifications, and optimize $A$ by minimizing the gradient variance
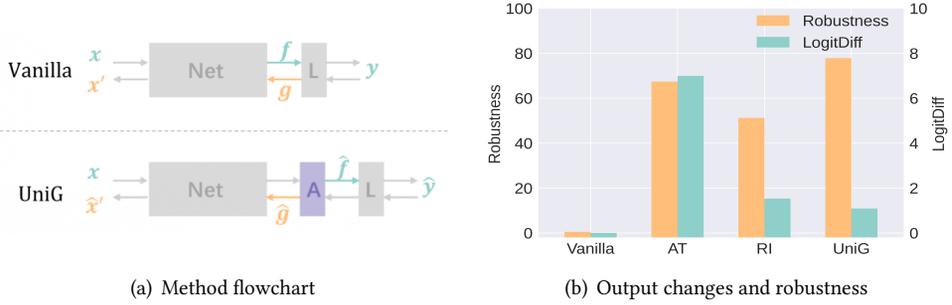
(a) Method flowchart



(b) Output changes and robustness

Fig. 1. (a) The flowchart of our UniG method. The plug-in Hadamard product module $A$ slightly modifies the forward feature ($\hat{f} \approx f$) but totally distorts the backward gradient ($\hat{g} \neq g$), where $\hat{g}$ is a less threatening attack direction. (b) The robust accuracy (under a SOTA SQA named Square [3]) and the logit (output) difference on clean data. Our UniG has both the best robustness and smallest output modifications compared with SOTA AT [22] and RI [44] defenses.

between different images to transform $\hat{g}$ to be a UAP direction. As a result, our method conceals the vulnerable attack direction $g$ and instead displays a weaker one $\hat{g}$ through the output probabilities $\hat{y}$ affected by the optimized module $A$. The advantages of our approach are listed below.

- The module $A$ is plug-and-play for any pretrained network with negligible additional computation. Each time a test batch arrives, according to our experiments, the module parameters can be well re-optimized from random initialization within one epoch using the current test data (Noting that the test data is unlabeled, we compute the cross-entropy loss by using the predicted class as the label). UniG is remarkably lightweight compared to other dynamic defenses [51, 52] that optimize the entire model in inference time.
- In contrast to other defenses, our UniG performs best in improving real-world robustness with the least modification on output probabilities, see Figure 1(b). Our explicit optimization objective of changing the gradient information contained in the output fundamentally suppresses the attack performance of SQAs, while other defenses such as RI randomly change the forward output to indirectly perturb the attacker's estimation on the gradient.

We compare our UniG with seven defenses [2, 22, 25, 37, 40, 44, 46, 51] on CIFAR10 [30] and ImageNet [19] under five popular SQAs [1, 3, 23, 27, 28]. The result shows that our UniG can keep clean accuracy (94.26% on CIFAR10 [30] and 78.47% on ImageNet [19]) while obtaining the best robustness in most cases (the remaining accuracy under 2500-query Square attack is 77.80% on CIFAR10, while the state-of-the-art (SOTA) AT model only achieves 67.34%). Additionally, under adaptive attacks in black-box scenarios such as steal-based [29, 38] and hyperparameter-tuning SQAs [44], our UniG also protects DNNs best compared to other defenses [44, 51]. Note that UniG as a black-box defense, while modifying gradients, does not belong to obfuscated gradients defenses [4] which resist white-box attacks effectively but cannot mitigate black-box ones.

## 2  RELATED WORK

### 2.1  Score-based Query Attacks

Current SQAs can be categorized into two types: gradient estimation [1, 12, 27, 28] and random direction search [3, 23]. For the first type, attackers calculate directional derivation using the output of continuous queries to estimate the input gradient [12]. Based on this idea, Bandits [28] further reduces query times through utilizing data-dependent gradient prior information. Natural

Evolutionary Strategies (NES, [27]) considers the condition of limited information, eg. only top-k probabilities are accessible, and develops the corresponding SQA. SignHunter [1] constructively proposes to ignore gradient magnitudes and only focus on gradient signs to improve attack efficiency. For the second type, attackers randomly choose an attack direction and adjust or directly abandon it according to the feedback from models. Square [3] is the most popular approach in this category, which adds localized square-shape random noises into inputs. In addition to Sqaure, SimBA [23] is an early and simple random-search attack, which selects perturbations from orthonormal basis.

## 2.2 Adversarial Defense

Before, adversarial defenses are mainly designed to resist white-box attacks [9, 36, 48], among which adversarial training (AT [22, 49]) is the most popular and comprehensive. Recently, defenses for SQAs are proposed in quantity, which can be categorized into three types: adversarial detection, denoising and randomness injection. The first one utilizes the similarity between malicious query data to detect adversarial examples [13, 32, 39]. However, it takes large storage resources and meets problem when facing long-interval queries. The basic idea of denoising methods is to pre-process inputs to eliminate adversarial perturbations [2, 37, 40]. Consequently, outputs for continuous queries remain unaltered, so that SQA attackers cannot obtain valid information of gradients. For instance, [37] proposes to use a pre-trained diffusion model as the input denoising network. [2] adds inverse adversarial perturbations to inputs, and [40] mixups the input with other random clean samples, all for the purpose of shrinking the adversarial perturbation. In contrast, randomness injection (RI) approaches aim to bewilder attackers by random noises on the outputs. To achieve this, they inject random noises into different parts of models such as inputs [8, 44, 54], features [34] and parameters [25]. In addition, [17, 31, 47] propose to combine randomness into training process for certified robustness. Although randomness protects models from attacking, it inevitably reduces clean accuracy. Apart from the above defenses, recently proposed transductive methods [51, 52] are highly related to our method, which dynamically optimize network parameters at test time to adapt attacks. Nevertheless, the test speed of dynamic defenses is slow because of their optimization process at each inference time. Different from them, our method only needs to adjust the parameter of our designed module instead of the whole network, thus is more practical. Although our approach, like the others, defends against SQAs by changing the output probabilities, the output modification in our method is obtained via optimizing a designed goal of keeping forward information and concealing backward knowledge, instead of random noises. This is why we obtain better accuracy and defense performance compared to RI methods which do not have a clear optimization goal.

## 3 METHOD

### 3.1 Preliminaries and Motivation

Let us denote the victim model as $M : X \rightarrow Y$, and the benign data as $(x, y) \in (X, Y)$. An attack aims to craft an adversarial example $x'$ which locates at the neighborhood of $x$ but misguides the model to an incorrect prediction class. The attack optimization problem can be summarized as

$$\min_{x' \in N_r(x)} l(x') = \min_{x' \in N_r(x)} (M_y(x') - \max_{j \neq y} M_j(x')), \tag{1}$$

where $N_r(x) = \{x' | \|x - x'\|_p \leq r\}$ indicates the $l_p$ ball around $x$ with a radius $r$, and $M_y$ denotes the predicted logits (or probabilities) of the $y$-th class. An effective attack algorithm can make $l(x') \leq 0$ so that the predicted class is no longer the true label $y$. A common method to solve (1) is to iteratively optimize the objective function by gradients of $l(x')$ w.r.t. $x'$, namely projected gradient descent (PGD, [36]) algorithm.

However, for SQAs, the gradient cannot be obtained directly since only the output probability of the victim model is accessible to attackers. Therefore, SQAs use direction derivation or random search method to estimate the gradient. The key of these methods is to utilize the forward output of queries to infer backward gradient information. Accordingly, we propose a defense idea that modifies the forward outputs slightly to prevent attackers from estimating backward gradients.

## 3.2 UniG: Unifying Data's Gradients to Defend against SQAs

We intend to change the gradient information contained in the output of queries to fool the attacker into a distorted attack trajectory, and simultaneously keep the output of clean data as far as possible. The overall goal of our method can be summarized as follows:

$$
\begin{cases} \hat{M}(x) \approx M(x) \\ \hat{G}(x) \neq G(x) \end{cases} \tag{2}
$$

where $M(x)$, $\hat{M}(x)$ denote the output logits (probabilities) of the vanilla and defense model respectively, and $G(x)$, $\hat{G}(x)$ are the backward gradient of inputs of the vanilla and defense model. We use the gradient constraint to guide our slight modification on outputs. The gradient information contained in our slightly modified output $\hat{M}(x)$ is the artful direction $\hat{G}(x)$ instead of the true vulnerable trajectory $G(x)$. The SQA attacker hereby can only dig an attack direction $\hat{G}(x)$ from the output of queries to generate AEs, which is designed to be weaker than $G(x)$. Here, we choose $\hat{G}(x)$ to be the direction of the universal attack perturbation [58] which is consistent between different samples and proved to be weaker [5, 56–58] than normal attack directions which are image-specific [36]. According to the above discussion, we propose an optimization problem to generate our modification on outputs to defend against SQAs.

$$
\begin{aligned}
\min_{\hat{M}} \quad & \sum_{i=1}^{n} (\hat{G}(x_i) - \hat{G}(x_{i+1}))^2 \\
\text{s.t.} \quad & \|\hat{M}(x) - M(x)\| \leq \delta,
\end{aligned} \tag{3}
$$

The constraint obviously corresponds to the slight modification target, where $\delta$ controls the degree of output offset. The objective function corresponds to the goal of distorting the gradient information into a universal one which is less threatening, where $\hat{G}(x_i)$ represents the gradient of the $i$-th input.

To solve problem (3), a direct way is to finetune the victim model using the above loss function. However, this solution is not only resource-consuming, but also difficult because of millions of network parameters and training samples. Therefore, we propose an alternative solution. We replace $\hat{M}$ with a simple module $A$, which is inserted into the penultimate layer of the victim model as shown in Figure 1(a). The parameter of the victim model remains unchanged and only the module $A$ needs optimization. Moreover, we choose to unify the gradient of features instead of inputs to simplify the optimization process. The overall problem becomes the following formulation,

$$
\begin{aligned}
\min_{A} \quad & l(A, x) = \sum_{i=1}^{b-1} (\hat{g}(x_i) - \hat{g}(x_{i+1}))^2 \\
\text{s.t.} \quad & \|\hat{f} - f\| \leq \delta.
\end{aligned} \tag{4}
$$

where $\hat{f}$, $\hat{g}$, and $b$ respectively denote feature, feature gradients, and batch size as shown in Figure 1(a). The operation of $A$ is designed to do Hadamard product with input features, that is, $\hat{f} = A \circ f$, which is simple but effective. To keep forward outputs, every element of $A$ is expected to be close to one, while the objective loss in problem (4) gives an instructional direction to optimize $A$ around

all-one matrix. As a result, the slight modification on features plays an important role in distorting backward information to mislead attackers.

We solve problem (4) using gradient descent (GD) algorithm, and the process of optimizing $A$ is integrated into the forward calculation of our defense model, see Alg. 1. At each inference time, the module $A$ is re-optimized from random Gaussian initialization using current test data with the objective function in Eq. (4). Noticing that test data have no labels, we utilize the prediction label as the true label in cross-entropy loss to calculate $\hat{g}(x)$. It is worth noticing that the objective loss is easy to optimize because only the linear layer and our designed module are involved in the calculation. Our experiment results show that with only one iteration step, the gradient variance can be efficiently minimized to significantly improve the robustness under SQAs.

---

**Algorithm 1** Forward Calculation of UniG Model (The symbols correspond to those in Figure 1(a))

---

**Input:** batch data: $x$; optimization iterations: $p$; learning rate: $\alpha$; constraint parameter: $\delta$
**Output:** model prediction: $\hat{y}$

1:  Initialize parameters $A$ with $A \sim \mathcal{N}(1, 0.5)$
2:  **for** $i = 1$ to $p$ **do**
3:      Compute $f = Net(x), y = L(f), \hat{y} = L(A(f)), c = OneHot(y), CE\_loss = -\sum c_i \cdot log(\hat{y}_i)$
4:      Calculate and min-max normalize the gradient $\hat{g}$
5:      Compute the objective loss $l(A, x)$ in problem (4)
6:      Update $A$ with $A \leftarrow A - \alpha \cdot \frac{\partial l(A,x)}{\partial A}$
7:      Clip $A$ with $\|A_{ij} - 1\|_\infty \le \delta$
8:  **end for**
9:  Compute final output $\hat{y} = L(A(f))$ with optimized $A$
10:  **return** $\hat{y}$

---

### 3.3 Discussion

To further elucidate our method, we display the training process and the final value of $A$ as showing in Figure 2(a),3(b). From Figure 2(a), we observe that the gradient unification loss is well-optimized with one iteration step (The forward consistency loss is calculated by $\|\hat{f} - f\|$ to show the modification on forward features is tiny.). Figure 3(b) shows partial final values of the matrix $A$ (the parameter of the designed module), which demonstrates that although randomly initialized, the value of $A$ is highly dependent on the current batch data after optimization, thus is different for each query. The difference is because that gradients are naturally diverse for different data, so unifying them requires divergent values of the matrix $A$.

Since our method claims that we fool attackers into a universal attack direction, we visualize the AEs and adversarial perturbations of our UniG model to verify it. The result in Figure 3(a) directly illustrates our statement, where the adversarial noise of our UniG model keeps consistent for different images, while that of the vanilla network obviously is different for diverse images. Because of the distortion of image-dependent attack directions, in our defense, the marginal loss in Eq. (1) hardly decreases as the number of queries increases, see Figure 2(b).

Considering a possible situation where the model receives a single image to test, our approach, like other dynamic defenses based on batch optimization, needs a solution to deal with this condition. One possible way is that we could cascade several training data with the test sample to perform such optimization for the single test sample situation. We specifically discuss it with experiment results in Section 4.6.

(a) Objective function value *w.r.t* iteration step
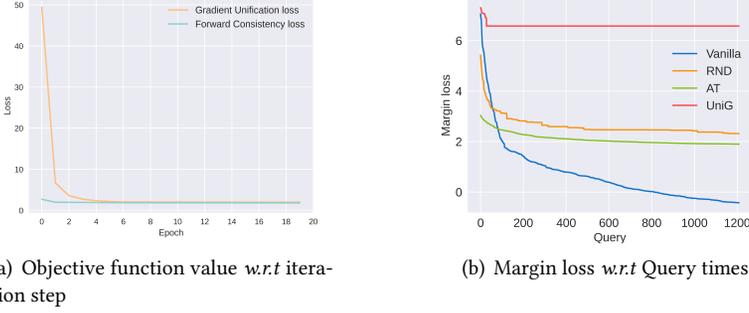
(b) Margin loss *w.r.t* Query times

Fig. 2. (a) Objective function value (UniG Loss in Eq. (4)) *w.r.t* iteration step of our defense model. Since each forward calculation the loss is re-optimized with the current test data, we randomly choose one forward process to present the loss change. (b) Margin loss *w.r.t* Query times of four compared defenses under Square attack. A higher loss means a more robust network.



(a) Adversarial examples and noises of UniG model

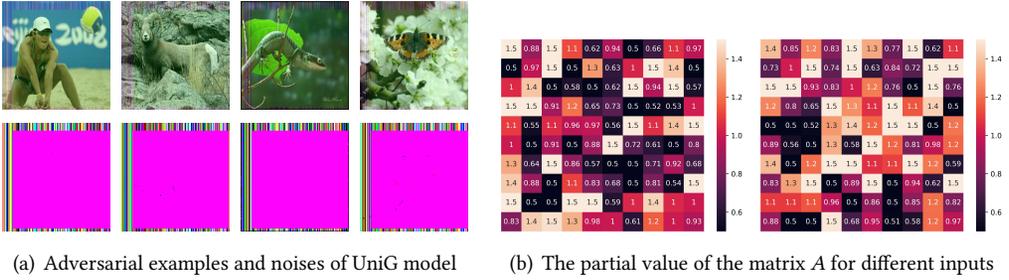(b) The partial value of the matrix $A$ for different inputs

Fig. 3. (a) Adversarial perturbations and the corresponding AEs of our UniG model for different images in one test batch, where the perturbations are the same for diverse images. (b) The partial value of the matrix $A$ for different inputs (queries), indicating that different data lead to divergent values of the matrix $A$.

## 4 EXPERIMENT

### 4.1 Setup

**Datasets.** We experiment on two widely-used classification datasets: CIFAR10 [30] and ImageNet [19]. For CIFAR10, the whole testset is used to evaluate method performances. For ImageNet, we randomly choose 1K images (one image for one class) from the validation set as test data.

**Models.** PreResNet18 [24] and WideResNet-50-2 [55] are used as typical model architectures. On CIFAR10, PreResNet18 with vanilla training achieves 94.26% accuracy. On ImageNet, we directely use the pre-trained model from torchvision package of PyTorch [42], which has 78.47% accuracy.

**Score-based Query Attacks.** We evaluate defense methods under five popular SQAs, including Square [3], SimBA [23], Sign [1], NES [27], and Bandits [28]. The first two are based on random search and the rest are based on gradient estimation. Most of the results in this paper are obtained under the untargeted $\ell_\infty$ norm attack, but we also present the results of the targeted and $\ell_2$ norm attack in Section Performance under More Attacks for completeness. The $\ell_\infty$ attack bound is set as 8/255 for CIFAR10 and 4/255 for ImageNet, and the $l_2$ attack bound is set as 1 for CIFAR10 and 5

for ImageNet. The query budget of SQAs is set to be 100 and 2500 for evaluating defense methods under different attack intensities.

**UniG setting.** The UniG module is plugged into vanilla training models. For its parameters, we set $\delta$=0.5, $p$=1, $\alpha$=10 on CIFAR10 and $\delta$=0.1, $p$=1, $\alpha$=1 on ImageNet.

**Compared defenses.** The proposed defense method will be compared with adversarial training (AT, [22, 46]), random noise defense (RND, [44]) and its enhancement (RND-GF), parameter noise injection (PNI, [25]), dynamic inference (DENT, [51]), mixup inference (Mixup, [40]), Anti-adversary combination (Anti-adv, [2]) and adversarial purification based on Diffusion models (DDPM, [37]) approaches. The first one is the most popular defense, the second to fourth ones belong to RI, the fifth one is a dynamic defense, and the rest are denoising methods. The AT models are obtained from Robustbench[1] [18] and we choose the most robust one for comparison. Random noise defense, which adds random noise into inputs or outputs, is conducted with noise variance 0.02, as recommended by [44]. RND-GF is fulfilled by fine-tuning the baseline model with 100 epochs, and then we test it using variance 0.05. The model of PNI is obtained from the corresponding github code[2] and we choose the best for comparison. DENT is a transductive method and we insert it into the vanilla training model, as the same as UniG. We try different mixup ratios and choose the best one 0.9 to combine Mixup[3] method with our baseline and AT model. For Anti-adv[4], we use the recommended iteration number $K = 2$ and the anti-adversary step is set to 4/255 for the best performance since the common test attack step is 8/255. We directly utilize the pre-trained diffusion model in DDPM code[5] to denoise inputs for each query.

**Metric.** For a good defense, we need to consider three folds: the clean accuracy, the logit difference, and the robust accuracy, *i.e.*, the remaining accuracy under attacks. The logit difference reflects the output difference between vanilla model and the current defense model, whose value is highly proportional to the probability difference and more pronounced to observe. In this paper, we use $l_2$ norm to measure it. Other norms can be used as well and the conclusion is similar. The smaller it is, the more friendly the model is to users.

## 4.2 Defense Performance

Table 1 comprehensively reports the defense performance of the proposed UniG together with AT, RND, RND-GF, PNI, DENT, Mixup, Anti-adv and DDPM. Due to the huge computing overhead of DDPM, we just test its performance under the Square attack which is classical and popular. AT could improve the robust accuracy under different attacks but it degrades clean accuracy by about 7% for CIFAR10 and 10% for ImageNet. The other defenses undertake effort to avoid the downgrade on clean accuracy, but the robust accuracy is generally lower than AT, especially on CIFAR10. Although DENT performs well on ImageNet under various SQAs, it sacrifices the accuracy of output probabilities. In fact, we discover that the output probability of DENT model is almost composed of zeros or ones, making SQAs degrade into DQAs and hereby achieving good performance. Another defense named Mixup also outperforms our approach in some cases on ImageNet, however, it requires thirty additional forward calculations to perform its stochastic input mixing and output ensemble averaging, which adds too much computation burden in practice. By contrast, our UniG model achieves similar defense performance as AT, actually for most cases UniG is more robust than AT, and meanwhile greatly remains the clean accuracy, almost the same as the vanilla training. For logits difference, the advantage of UniG is more significant that the logits difference is generally

---

[1]https://robustbench.github.io/
[2]https://github.com/elliothe/CVPR_2019_PNI
[3]https://github.com/P2333/Mixup-Inference
[4]https://github.com/MotasemAlfarra/Combating-Adversaries-with-Anti-Adversaries
[5]https://github.com/NVlabs/DiffPure

Table 1. The comparison of different defense methods under $l_\infty$ norm attacks (query = 100/2500) on CIFAR10 and ImageNet. The clean accuracy, logits difference and robust accuracy are reported. The higher the clean accuracy and robust accuracy, the smaller the logits difference, indicating the better performance of the defense method. The best results are in bold, and the 2nd ones are underline.

| Datasets | Methods | Clean | Logit-diff | Square | SimBA | Sign | NES | Bandits |
|---|---|---|---|---|---|---|---|---|
| CIFAR10 | Vanilla | 94.26 | - | 38.79/0.46 | 41.03/0.45 | 48.25/0.26 | 75.28/10.06 | 68.81/25.92 |
| | AT | 87.35 | 7.00 | 79.15/67.34 | 83.61/71.36 | 78.28/64.43 | 85.30/79.49 | 83.90/74.33 |
| | RND | 91.14 | 1.53 | 65.04/51.22 | 74.88/63.07 | 64.71/51.95 | 85.67/69.27 | 67.79/58.27 |
| | RND-GF | 92.87 | 1.81 | 78.02/69.08 | 82.67/75.45 | 72.44/63.15 | 89.16/81.73 | 82.07/74.89 |
| | PNI | 85.93 | 8.16 | 64.66/51.54 | 64.13/57.41 | 65.31/50.48 | 80.77/69.30 | 66.93/56.81 |
| | DENT | 94.25 | 7.35 | 81.78/57.71 | 75.23/54.00 | 64.09/46.18 | 88.60/67.80 | 74.55/69.75 |
| | Mixup | 94.32 | 8.43 | 76.97/37.72 | 75.50/54.82 | 51.88/15.09 | 90.55/75.46 | 78.91/73.96 |
| | Anti-adv | 92.63 | 3.23 | 62.45/30.31 | 57.70/31.54 | 52.80/29.64 | 86.14/58.36 | 70.55/66.05 |
| | DDPM | 88.88 | 2.30 | 52.50/42.80 | - | - | - | - |
| | UniG | 94.26 | 1.09 | 81.90/77.80 | 89.79/86.42 | 72.58/68.81 | 89.55/67.88 | 80.31/75.86 |
| ImageNet | Vanilla | 78.47 | - | 52.71/6.70 | 64.66/6.28 | 51.00/11.77 | 68.27/59.64 | 67.64/35.23 |
| | AT | 68.41 | 49.23 | 62.25/51.92 | 67.50/51.92 | 60.88/56.10 | 65.67/65.67 | 66.49/61.64 |
| | RND | 77.14 | 13.73 | 61.13/48.44 | 71.22/65.25 | 59.40/54.77 | 74.83/72.51 | 70.02/67.42 |
| | DENT | 78.60 | 98.73 | 68.85/63.47 | 77.69/77.16 | 70.74/59.74 | 77.03/74.67 | 76.62/74.25 |
| | Mixup | 77.95 | 62.06 | 69.75/54.57 | 77.65/- | 68.60/- | 76.39/- | 75.44/- |
| | Anti-adv | 72.92 | 49.76 | 50.04/28.40 | 62.04/53.27 | 53.23/42.29 | 69.27/61.40 | 63.44/61.89 |
| | UniG | 78.47 | 2.75 | 66.14/52.88 | 78.22/77.44 | 57.28/45.51 | 77.68/71.40 | 74.14/72.65 |

reduced by an order of magnitude, which means UniG is much more user-friendly than the other defenses.

As a plug-and-play module, UniG can be readily applied in any network, for example in a model trained by adversarial examples. Besides, since Mixup approach is based on AT models, we also conduct it to compare our performance here. The performance of UniG and Mixup in AT can be found in Table 2. From AT, the classification accuracy of UniG-AT is well kept on clean examples and only drops 3.04% on CIFAR10 under 2500-query Square attack, while the AT model drops about 20%. Similar performance could be found on both CIFAR10 and ImageNet under different attacks. Besides, UniG surpasses Mixup in terms of robustness both on CIFAR10 and ImageNet. The improvement of the Mixup method is unstable, seeing the drop in robust accuracy on ImageNet.

Table 2. The performance of AT-based UniG model under SQAs (query = 100/2500) on CIFAR10 and ImageNet.

| Datasets | Methods | Clean | Square | SimBA | Sign | NES | Bandits |
|---|---|---|---|---|---|---|---|
| CIFAR10 | AT | 87.35 | 79.15/67.34 | 83.61/71.36 | 78.28/64.43 | 85.30/79.49 | 83.90/74.33 |
| | Mixup-AT | 86.53 | 84.91/82.56 | 86.59/86.29 | 79.61/77.88 | 85.66/84.80 | 84.71/83.99 |
| | UniG-AT | 87.36 | 84.61/84.32 | 86.98/86.98 | 83.87/83.87 | 87.36/86.49 | 75.47/73.03 |
| ImageNet | AT | 68.41 | 62.25/51.92 | 67.50/51.92 | 60.88/56.10 | 65.67/65.67 | 66.49/61.64 |
| | Mixup-AT | 67.79 | 59.26/52.20 | 60.54/- | 64.40/- | 67.40/- | 67.34/- |
| | UniG-AT | 68.36 | 65.49/63.71 | 68.29/68.09 | 64.94/64.94 | 68.36/67.68 | 67.60/66.92 |

## 4.3 Performance under More Attacks

Before, we have shown the performance under current popular SQAs in an untargeted setting and a $\ell_\infty$-norm bound. In this section, we further evaluate the performance of our method under

other settings, *i.e.*, $\ell_2$ / $\ell_\infty$ norm and target (-T) / untarget (-UT) attacks as listed in the first row of Table 3. We adopt Square attack, the current SOTA SQA attack, for robustness evaluation. As Table 3 verified, plugging UniG into vanilla and AT model can both achieve significant robustness improvement with a little clean accuracy drop. For instance, on CIFAR10, our UniG-Vanilla model improves robust accuracy by no less than 71.29% under all attack settings with 2500 queries, and UniG-AT improves by $\geq$ 13.99% compared with AT model, while our method keeps the same or even higher clean accuracy.

Table 3. UniG under $\ell_2/l_\infty$ norm and target (-T) /untarget (-UT) Square attack (query = 100/2500)

| Datasets | Methods | Clean | $\ell_\infty$-T | $\ell_\infty$-UT | $\ell_2$-T | $\ell_2$-UT |
|---|---|---|---|---|---|---|
| CIFAR10 | Vanilla | 94.26 | 45.24/0.35 | 38.79/0.46 | 56.56/6.60 | 62.21/5.66 |
| | UniG | **94.26** | **78.26/71.64** | **81.90/77.80** | **85.78/84.83** | **87.66/80.12** |
| | AT | 87.35 | 77.74/66.37 | 79.15/67.34 | 76.87/65.51 | 81.23/68.13 |
| | UniG-AT | **87.36** | **83.87/83.87** | **84.61/84.32** | **79.50/79.50** | **84.74/84.74** |
| ImageNet | Vanilla | 78.47 | 59.79/14.91 | 52.71/6.70 | 50.30/7.14 | 41.59/7.30 |
| | UniG | **78.47** | **64.35/51.79** | **66.14/52.88** | **55.95/38.21** | **56.97/40.41** |
| | AT | **68.41** | 62.25/53.77 | 62.25/51.92 | 55.14/41.63 | 55.07/39.61 |
| | UniG-AT | 68.36 | **65.35/65.08** | **65.49/63.71** | **62.21/59.74** | **61.66/59.68** |

## 4.4 Adaptive Attack

In general, an adaptive attack refers to an elaborate white-box attack with full knowledge of the defense strategy. Nevertheless, in the real-world cases, the attacker and the victim model are double-blind to each other, which means the defense strategy is not prior information for the attacker. Therefore, we consider another two adaptive attacks for black-box conditions: model stealing for transferable attacks [29, 38] and tuning hyper-parameters for optimal attacks [44].

The first one could utilize model outputs to estimate a surrogate model and then use the white-box adversarial examples of the surrogate model to attack the original model based on attack transferability [41]. We here adopt two classical and practical steal-based attacks, eg. MAZE[6] [29] and KnockoffNet[7] [38], to verify the effectiveness of our method. We use the recommended hyper-parameters in [29], and the query budget is set to $3 \times 10^7$. Table 4 demonstrates that UniG still outperforms baseline and RND by > 20% under strong KnockoffNet attack. Although DENT outperforms UniG under the KnockoffNet attack, its performance under the more practical MAZE attack, which does not require a surrogate dataset as KnockoffNet does, is even worse than the baseline model, while our UniG improves the robust accuracy by 20.2%.

For the optimal hyper-parameter attack, we follow the settings in RND [44] and choose different hyper-parameters of Square, NES and Bandits attacks to find the optimal one. The result is shown in Table 5, 6. As the square size or update step increases, our defense robustness decreases slightly, but within the normal fluctuation range as demonstrated by the performance of RND.

---

[6]https://github.com/sanjaykariyappa/MAZE
[7]https://github.com/tribhuvanesh/knockoffnets

Table 4. Remaining accuracy of RND, DENT, and UniG methods on CIFAR10 under steal-based adaptive attacks.

|  | Clean | MAZE | KnockoffNets |
|---|---|---|---|
| Baseline | 94.26 | 65.00 | 21.95 |
| RND | 90.91 | 84.05 | 30.02 |
| DENT | 94.25 | 59.06 | **64.14** |
| UniG | **94.26** | **85.20** | 50.28 |

Table 5. Remaining accuracy of RND and UniG *w.r.t.* square size of Square attack (query=100/2500) on CIFAR10.

| square size | RND | UniG |
|---|---|---|
| 0.05 | 65.04/51.22 | **81.90/77.80** |
| 0.1 | 63.25/47.64 | **80.99/74.89** |
| 0.2 | 62.01/42.89 | **79.64/72.16** |
| 0.3 | 60.83/41.05 | **78.79/70.11** |

Table 6. Remaining accuracy of RND / UniG *w.r.t.* update step of NES / Bandits (query=100/2500) on CIFAR10.

| Attack | update step | 0.0001 | 0.0005 | 0.001 | 0.005 | 0.01 | 0.05 |
|---|---|---|---|---|---|---|---|
| NES | RND | 87.49/79.29 | 86.58/78.38 | 86.58/78.38 | 85.79/**71.09** | 83.85/**68.36** | 74.73/59.24 |
|  | UniG | **94.26/91.43** | **94.26/86.72** | **93.32/82.95** | **90.49**/69.75 | **87.66**/66.92 | **79.18/65.98** |
| Bandits | RND | 79.78/61.07 | 77.85/58.90 | 74.09/58.63 | 68.77/58.19 | 67.39/57.83 | 65.36/58.15 |
|  | UniG | **94.03/86.15** | **89.45/82.34** | **83.78/78.79** | **80.34/75.95** | **80.34/75.94** | **78.48/76.47** |

## 4.5 Hyper-parameter Study

Previous experiments are conducted using the fixed hyper-parameters as introduced before. To further evaluate our method, we here study the influence of different hyper-parameters on the performance of our method. The main hyper-parameters include the forward constraint parameter $\delta$, the optimization iteration $p$, and the learning rate $\alpha$. The parameter $\delta$ decides the largest element-wise difference between the parameters of our module and all-one matrix, which controls the trade-off between the clean performance and the robustness of our model. The iteration $p$ and the learning rate $\alpha$ are related to the optimization process of our module $A$. For fast test speed, we adopt one iteration and relatively bigger $\alpha$ to optimize $A$ in the previous experiments, while in this section, we adopt different $p$ and $\alpha$ to observe their influence on clean accuracy, logit difference and robustness under Square attack. As Fig 4(a) shows, with the increase of $\alpha$, the logit difference slightly increases, but is still less than which of other methods. And for the optimization iteration $p$, we observe the same phenomenon as $\alpha$, see Fig 4(b). For the change of $\delta$, the robustness keeps stable, while the clean accuracy suffers from small fluctuations and the logit difference slightly increases, see Fig 4(c). Apart from the aforementioned hyper-parameters, we also evaluate our performance under different batch sizes, since our objective loss is dependent on the input test data. The result in Figure 4(d) reveals that our method performs better with a larger batch size, but even with the batch size equals to 32, our robustness performance is quite good.

## 4.6 Other Algorithmic Discussions

This section discusses the computational complexity and the single test sample situation of our method. The first one is related to the time and space overhead of our method. And the single test sample condition is necessary to study taking the reality into account.

**Computational Complexity.** Because of the optimization process contained in our model's forward calculation, the computational complexity of our approach inevitably increases. There are two metrics to measure this complexity — the number of model parameters and the number of Floating Point Operations (FLOPs) of the forward process. The first one reflects the space overhead and the second one represents the time overhead. Table 7 reports the computational complexity of
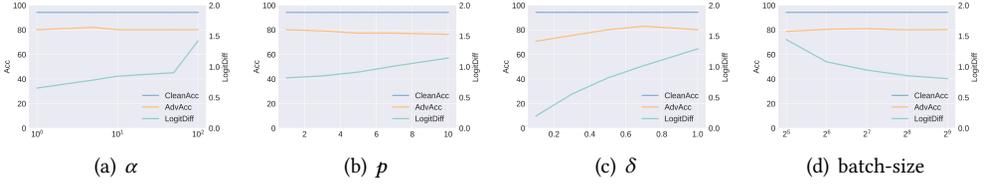
(a) $\alpha$         (b) $p$         (c) $\delta$         (d) batch-size

Fig. 4. The influence of hyper-parameters on our method. The main hyper-parameters of UniG are the forward constraint parameter $\delta$, the optimization iterations $p$, and learning rate $\alpha$, and the batch size $b$, which influences our objective loss in batch-wise optimization.

Table 7. The model complexity of vanilla, DENT and UniG network. High FLOPs indicate high complexity.

| Methods | Parameter (M) | FLOPs (GMac) |
|---------|---------------|--------------|
| Vanilla | 11.17 | 0.54 |
| DENT | 11.17 | 7.78 |
| UniG | 11.43 | 0.55 |

Vanilla, DENT and UniG model on CIFAR10 with the PreResNet18 architecture. Although there is some growth in the number of network parameters, the FLOPs of our method is far less than that of DENT, with a slight increase compared with that of the vanilla model.

**Single Test Sample.** Since it is possible that some times the network receives a single image, it is imperative to discuss our performance in this case. Although our UniG, as well as most transduction defenses like DENT, are based on batch optimization, we propose a solution that cascading several (we use ten in our experiments) training data with test data to perform such optimization for the single test sample situation. Based on the experiment result on CIFAR10 in Table 8, we can conclude that although there is a debasement of robustness when batch size equals to one, our method still surpasses baseline and DENT, and even slightly outperforms RND. We use Square attack and set $p = 0.05$, query budget=100 here. Thus, although the primary application of our approach is not for the single test sample situation, our method can still improve robustness for this case.

Table 8. Remaining accuracy of DENT and UniG under Square attack (query=100) on CIFAR10 for the case of a single test sample. The last two columns are used to compare the performance.

|        | Baseline | DENT(bz=1) | UniG(bz=1) | RND | UniG(bz=256) |
|--------|----------|------------|------------|-----|--------------|
| Clean  | 94.26 | 83.22 | **94.26** | 91.14 | **94.26** |
| Square | 38.79 | 5.32  | **72.24** | 65.04 | **81.90** |

## 5 CONCLUSION

In this paper, we propose a new defense method named Unifying Gradients (UniG) to defend against the most threatening attack in real applications—score-based query attacks (SQAs). The proposed method is based on the idea of distorting the gradient information contained in the query output by a slight modification on the forward output to fool the attacker into a weaker attack trajectory. In this paper, we choose the universal attack perturbation (UAP) as the weaker direction. Accordingly, the change on outputs is explicitly optimized with the gradient unification

loss which indicates the UAP path. To practically implement this modification, we propose a Hadamard product operation module, which can be inserted into any pre-trained networks, and optimizes its parameter with the designed forward consistency and backward distortion loss at each inference time. With comprehensive experiments on CIFAR10 and ImageNet, it is verified that our approach can significantly boost the robustness under SQAs with no sacrifice of clean accuracy and a few variation on clean outputs. Noticing that the designed module is plug-and-play with negligible extra computational overhead, the overall method has a promising application prospect in the future.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Abdullah Al-Dujaili and Una-May O'Reilly. 2020. Sign Bits Are All You Need for Black-Box Attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

[2] Motasem Alfarra, Juan C Pérez, Ali Thabet, Adel Bibi, Philip HS Torr, and Bernard Ghanem. 2022. Combating adversaries with anti-adversaries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 5992–6000.

[3] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. 2020. Square attack: a query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*. Springer, 484–501.

[4] Anish Athalye, Nicholas Carlini, and David Wagner. 2018. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*. PMLR, 274–283.

[5] Philipp Benz, Chaoning Zhang, Tooba Imtiaz, and In So Kweon. 2020. Double targeted universal adversarial perturbations. In *Proceedings of the Asian Conference on Computer Vision*.

[6] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. 2016. End to end learning for self-driving cars.

[7] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2018. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

[8] Junyoung Byun, Hyojun Go, and Changick Kim. 2022. On the Effectiveness of Small Input Noise for Defending Against Query-based Black-Box Attacks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3051–3060.

[9] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *the IEEE Symposium on Security and Privacy (SP)*. 39–57.

[10] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*. ACM, 1721–1730.

[11] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. 2020. HopSkipJumpAttack: A Query-Efficient Decision-Based Attack. In *the IEEE Symposium on Security and Privacy (SP)*. 1277–1294.

[12] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *the 10th ACM workshop on artificial intelligence and security*. 15–26.

[13] Steven Chen, Nicholas Carlini, and David Wagner. 2020. Stateful detection of black-box adversarial attacks. In *Proceedings of the 1st ACM Workshop on Security and Privacy on Artificial Intelligence*. 30–39.

[14] Sizhe Chen, Fan He, Xiaolin Huang, and Kun Zhang. 2022. Relevance attack on detectors. In *Pattern Recognition*. 108491.

[15] Sizhe Chen, Zhengbao He, Chengjin Sun, and Xiaolin Huang. 2022. Universal Adversarial Attack on Attention and the Resulting Dataset DAmageNet. In *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*. 2188–2197.

[16] Sizhe Chen, Zhehao Huang, Qinghua Tao, and Xiaolin Huang. 2021. QueryNet: Attack by Multi-Identity Surrogates. *arXiv preprint arXiv:2105.15010* (2021).

[17] Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. 2019. Certified Adversarial Robustness via Randomized Smoothing. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*. PMLR, 1310–1320.

[18] Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. 2021. RobustBench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

[19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA*. IEEE Computer Society, 248–255.

[20] Kun Fang, Qinghua Tao, Yingwen Wu, Tao Li, Jia Cai, Feipeng Cai, Xiaolin Huang, and Jie Yang. 2020. Towards Robust Neural Networks via Orthogonal Diversity.

[21] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

[22] Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy A. Mann. 2021. Improving Robustness using Generated Data. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*. 4218–4233.

[23] Chuan Guo, Jacob R. Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Q. Weinberger. 2019. Simple Black-box Adversarial Attacks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, Vol. 97. PMLR, 2484–2493.

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity mappings in deep residual networks. In *European conference on computer vision*. Springer, 630–645.

[25] Zhezhi He, Adnan Siraj Rakin, and Deliang Fan. 2019. Parametric Noise Injection: Trainable Randomness to Improve Deep Neural Network Robustness Against Adversarial Attack. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 588–597.

[26] Zhichao Huang and Tong Zhang. 2020. Black-Box Adversarial Attack with Transferable Model-based Embedding. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

[27] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. 2018. Black-box Adversarial Attacks with Limited Queries and Information. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. PMLR, 2142–2151.

[28] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. 2019. Prior Convictions: Black-box Adversarial Attacks with Bandits and Priors. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

[29] Sanjay Kariyappa, Atul Prakash, and Moinuddin K Qureshi. 2021. Maze: Data-free model stealing attack using zeroth-order gradient estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13814–13823.

[30] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).

[31] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. 2019. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 656–672.

[32] Huiying Li, Shawn Shan, Emily Wenger, Jiayun Zhang, Haitao Zheng, and Ben Y Zhao. 2020. Blacklight: Defending black-box adversarial attacks on deep neural networks.

[33] Tao Li, Yingwen Wu, Sizhe Chen, Kun Fang, and Xiaolin Huang. 2022. Subspace Adversarial Training. In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 13409–13418.

[34] Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. 2018. Towards robust neural networks via random self-ensemble. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 369–385.

[35] Xuanqing Liu, Yao Li, Chongruo Wu, and Cho-Jui Hsieh. 2019. Adv-BNN: Improved Adversarial Defense through Robust Bayesian Neural Network. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

[36] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

[37] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar. 2022. Diffusion Models for Adversarial Purification. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, Vol. 162. PMLR, 16805–16827.

[38] Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. 2019. Knockoff nets: Stealing functionality of black-box models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4954–4963.

[39] Ren Pang, Xinyang Zhang, Shouling Ji, Xiapu Luo, and Ting Wang. 2020. AdvMind: Inferring Adversary Intent of Black-Box Attacks. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*. ACM, 1899–1907.

[40] Tianyu Pang, Kun Xu, and Jun Zhu. 2020. Mixup Inference: Better Exploiting Mixup to Defend Adversarial Attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

[41] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. 506–519.

[42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. 8024–8035.

[43] Rafael Pinot, Laurent Meunier, Alexandre Araujo, Hisashi Kashima, Florian Yger, Cédric Gouy-Pailler, and Jamal Atif. 2019. Theoretical evidence for adversarial robustness through randomization. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. 11838–11848.

[44] Zeyu Qin, Yanbo Fan, Hongyuan Zha, and Baoyuan Wu. 2021. Random Noise Defense Against Query-Based Black-Box Attacks. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*. 7650–7663.

[45] Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C Duchi, and Percy Liang. 2019. Adversarial training can hurt generalization. *arXiv preprint arXiv:1906.06032* (2019).

[46] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. 2020. Do Adversarially Robust ImageNet Models Transfer Better?. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

[47] Hadi Salman, Jerry Li, Ilya P. Razenshteyn, Pengchuan Zhang, Huan Zhang, Sébastien Bubeck, and Greg Yang. 2019. Provably Robust Deep Learning via Adversarially Trained Smoothed Classifiers. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. 11289–11300.

[48] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.

[49] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. 2018. Ensemble Adversarial Training: Attacks and Defenses. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

[50] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. 2019. Robustness May Be at Odds with Accuracy. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

[51] Dequan Wang, An Ju, Evan Shelhamer, David Wagner, and Trevor Darrell. 2021. Fighting Gradients with Gradients: Dynamic Defenses against Adversarial Attacks.

[52] Yi-Hsuan Wu, Chia-Hung Yuan, and Shan-Hung Wu. 2020. Adversarial Robustness via Runtime Masking and Cleansing. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, Vol. 119. PMLR, 10399–10409.

[53] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. 2020. Adversarial examples improve image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 819–828.

[54] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan L. Yuille. 2018. Mitigating Adversarial Effects Through Randomization. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

[55] Sergey Zagoruyko and Nikos Komodakis. 2016. Wide Residual Networks. In *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*. BMVA Press.

[56] Chaoning Zhang, Philipp Benz, Tooba Imtiaz, and In-So Kweon. 2020. Cd-uap: Class discriminative universal adversarial perturbation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 6754–6761.

[57] Chaoning Zhang, Philipp Benz, Adil Karjauv, and In So Kweon. 2021. Data-free universal adversarial perturbation and black-box attack. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7868–7877.

[58] Chaoning Zhang, Philipp Benz, Chenguo Lin, Adil Karjauv, Jing Wu, and In So Kweon. 2021. A Survey on Universal Adversarial Attack. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*. International Joint Conferences on Artificial Intelligence Organization, 4687–4694.

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

[60] Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. 2020. Adversarial attacks on deep-learning models in natural language processing: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)* 11, 3 (2020), 1–41.