# Towards Visual Distortion in Black-Box Attacks

**Nannan Li and Zhenzhong Chen**[*]

School of Remote Sensing and Information Engineering, Wuhan University

## Abstract

Constructing adversarial examples in a black-box threat model injures the original images by introducing visual distortion. In this paper, we propose a novel black-box attack approach that can directly minimize the induced distortion by learning the noise distribution of the adversarial example, assuming only loss-oracle access to the black-box network. The quantified visual distortion, which measures the perceptual distance between the adversarial example and the original image, is introduced in our loss whilst the gradient of the corresponding non-differentiable loss function is approximated by sampling noise from the learned noise distribution. We validate the effectiveness of our attack on ImageNet. Our attack results in much lower distortion when compared to the state-of-the-art black-box attacks and achieves 100% success rate on InceptionV3, ResNet50 and VGG16bn. The code is available at https://github.com/Alina-1997/visual-distortion-in-attack.

## 1 Introduction

Adversarial attack has been a well-recognized threat to existing Deep Neural Network (DNN) based applications. It injects small amount of noise to a sample (e.g., image, speech, language) but degrades the model performance drastically [1, 2, 3]. With the continuous improvements of DNN, such attack could cause serious consequences in practical conditions where DNN is used. According to [4, 5], adversarial attack has been a practical concern in real-world problems, ranging from cell-phone camera attack to attacking self-driving cars.

According to the information that an adversary has of the target network, existing attack roughly falls into two categories: white-box attack that knows all the parameters of the target network, and black-box attack that has limited access to the target network. Each category can be further divided into several subcategories depending on the adversarial strength [6]. The proposed attack in this paper belongs to loss-oracle based black-box attack, where the adversary can obtain the output loss from supplied inputs. In real-world scenario, it's sometimes difficult or even impossible to have full access to certain networks, which makes the black-box attack practical and attract more and more attention.

Black-box attack has very limited or no information of the target network and thus is more challenging to perform. In the $l_p$-bounded setting, a black-box attack is usually evaluated on two aspects: number of queries and success rate. In addition, recent work [7] shows that visual distortion in the adversarial examples is also an important criteria in practice. Even under a small $l_\infty$ bound, perturbing pixels in the image without considering the visual impact could make the distorted image very annoying. As shown in Fig. 1, an attack [8] under a small noise level ($l_\infty \leq 0.05$) causes relatively large visual distortion and the perturbed image is more distinguishable from the original one. Therefore, under the assumption that the visual distortion caused by the noise is related to the spatial distribution of the perturbed pixels, we take a different view from previous work and focus on *explicitly* learning a noise distribution based on its corresponding visual distortion.

In this paper, we propose a novel black-box attack that can directly minimize the induced visual distortion by learning the noise distribution of the adversarial example, assuming only loss-oracle access to the black-box network. The quantified visual distortion, which measures the perceptual distance between the adversarial example and the original image, is introduced in our loss where the gradient of the corresponding non-differentiable loss function is approximated by sampling noise from the learned noise distribution. The proposed attack can achieve a trade-off between visual distortion and query efficiency by introducing the weighted perceptual distance metric in addition to the original loss. Theoretically, we prove the convergence of our model under a convex or non-convex loss function. The experiments demonstrate the effectiveness of our attack on ImageNet. Our attack results in much lower distortion than the other attacks and achieves 100% success rate on InceptionV3, ResNet50 and VGG16bn. In addition, it is shown that our attack is valid even when it's only allowed to perturb pixels that are out of the target object in a given image.

Our contributions are as follows:

- We are the first to introduce perceptual loss in a *non-differentiable* way for the generation of less-distorted adversarial examples. And the proposed method can achieve a trade-off between visual distortion and query efficiency by using the weighted perceptual distance metric in addition to the original loss.
- Theoretically, we prove the convergence of our model.
- Through extensive experiments, we show that our attack results in much lower distortion than the other attacks.

## 2 Related Work

Recent research on adversarial attack [9, 10, 11] has made advanced progress in developing strong and computationally efficient adversaries. In the following, we briefly introduce existing attack techniques in both the white-box and black-box settings.

### 2.1 White-box Attack

In white-box attack, the adversary knows the details of a network, including network structure and its parameter values. Goodfellow *et al.* [12] proposed a fast gradient sign method to generate
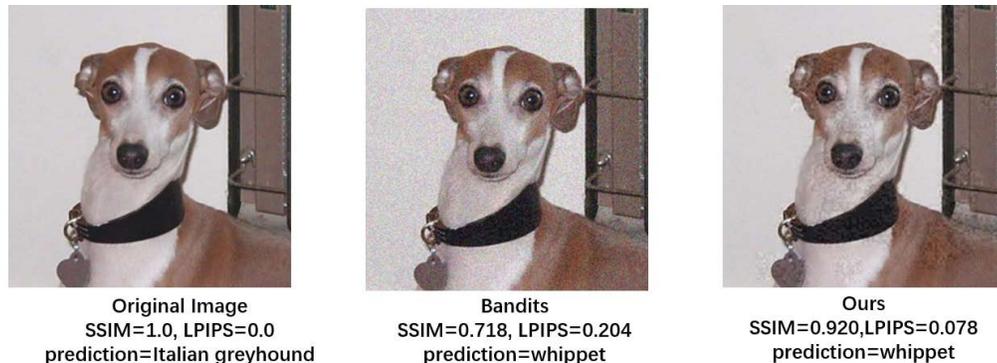
**Original Image**
SSIM=1.0, LPIPS=0.0
prediction=Italian greyhound

**Bandits**
SSIM=0.718, LPIPS=0.204
prediction=whippet

**Ours**
SSIM=0.920,LPIPS=0.078
prediction=whippet

Figure 1: Adversarial examples on ImageNet with bounded noise $\|\delta\|_\infty \leq 0.05$. The first image is the original unperturbed image. The following examples are from [8] and our method, respectively. Higher Structural SIMilarity (SSIM) and lower Learned Perceptual Image Patch Similarity (LPIPS) indicate less visual distortion.

adversarial examples. It's computationally effective and serves as a baseline for attacks with additive noise. In [13], a functional adversarial attack that applied functional noise instead of additive noise to the image, is introduced. Recently, Jordan *et al.* [7] stressed quantifying perceptual distortion of the adversarial examples by leveraging perceptual metrics to define an adversary. Different from our method which directly optimizes the metric, their model conducts a search over parameters from several composed attacks. There are also attacks that sample noise from a noise distribution [14, 15], on the condition that gradients from the white-box network is accessible. Specifically, [14] utilizes particle approximation to optimize a convex energy function. [15] formulates the attack problem as generating a sequence of adversarial examples in a Hamiltonian Monte Carlo framework.

In summary, white-box attack is hard to detect or defend [16]. In the meantime, however, it suffers from label-leaking and gradient-masking problem [2]. The former causes adversarially trained models to perform better on adversarial examples than original images, and the latter neutralizes the useful gradient for adversaries. The preliminary of acquiring full access to a network in white-box attack is also sometimes difficult to satisfy in real-world scenarios.

### 2.2 Black-box Attack

Black-box attack considers the target network as a black-box, and has limited access to the network. We discuss loss-oracle based attack here, where the adversary assumes only loss-oracle access to the black-box network.

**Query Efficient Attacks.** Attacks of this kind roughly fall into three categories: 1) Methods that estimate gradient of the black-box. Some methods estimate the gradient by sampling around a certain point, which formulates the task as a problem of continuous optimization. Tu *et al.* [17] searched for perturbations in the latent space of an auto-encoder. [18] utilizes feedback knowledge to alter the searching directions for efficient attack. Ilyas *et al.* [8] exploited prior information about the gradient. Al-Dujaili and O'Reilly [9] reduced query complexity by estimating just the sign of the gradient. In [19, 20], the proposed methods perform search in a constructed low-dimensional space. [21] shares similarity with our method as it also *explicitly* defines a noise distribution. However, the distribution in [21] is assumed to be an isometric normal distribution without considering visual distortion whilst our method does not assume the distribution to be a specific form. We compare with their method in details in the experiments. Other approaches in this category develop a substitute model [3, 22, 23] to approximate performance of the black-box. By exploiting the transferability of adversarial attack [12], the white-box attack technique applied to the substitute model can be transferred to the black-box. These approaches assume only label-oracle to the targeted network, whereas training of the substitute model requires either access to the training dataset of the black-box, or collection of a new dataset. 2) Methods based on discrete optimization. In [24, 9], an image is divided into regular grids and the attack is performed and refined on each grid. Meunier *et al.* [25] adopted the tiling trick by adding the same noise for small square tiles in the image. 3) Methods that leverage evolutionary strategies or random search [25, 26]. In [26], the noise value is updated by a square-shaped random search at each query. Meunier *et al.* [25] developed a set of attacks using evolutionary algorithms using both continuous and discrete optimization.

**Attacks that Consider Visual Impact.** Query efficient black-box attacks usually do not consider the visual impact of the induced noise, for which the adversarial example could suffer from significant visual distortion. Similar to our work, there are research that address the perceptual distance between the adversarial examples and the original image. [27, 28] introduce Generative Adversarial Network (GAN) based adversaries, where the gradient of the perceptual distance in the generator is computed through backpropagation. [29, 30] also require the adopted perceptual distance metric to be differentiable. Computing the gradients of a complex perceptual metric at each query might be computationally expensive [31], and is not possible for some rank-based metrics [32]. Different from these methods, our approach treats the perceptual distance metric as a black-box, saving the efforts of computing its gradients, and minimizing such distance by sampling from a learned noise distribution. On the other hand, [33, 34] present semantic perturbations for adversarial attacks. The produced noise map is semantically meaningful to human, whilst the image content of the adversarial example is distinct from that of the original

---

**Algorithm 1:** Our Algorithm

---
**Input:** image $x$, maximum norm $\epsilon$, proportion $q$ of the resampled noise
**Output:** adversarial example $x + \delta$

1   Initialize noise distribution $p_{\theta_0} = \text{softmax}(\theta_0)$ and noise $\delta_0$
2   **for** step $t$ in $\{1, ..., n\}$ **do**
3      $T^* = \text{argmin}_{T=0,1,...t-1} L(x, x + \delta_T)$
4      Compute baseline $b = L(x, x + \delta_{T^*})$
5      Update $\theta$ using Eq. (14), $\theta_t \leftarrow \theta_{t-1} - \nabla F(\theta_{t-1})$
6      Sample $\delta_t$, $\delta_t \leftarrow \text{resample}(\delta_{T^*}, q; \delta_{t-1})_{\delta_{t-1} \sim p_{\theta_{t-1}}}$
7      **if** $successful\_attack(x, x + \delta_t)$ **then**
8         return $x + \delta_t$

9   **def** $successful\_attack(x, x + \delta_t)$**:**
10      **if** $argmax_{k_1} f(x + \delta_t)_{k_1} \neq argmax_{k_2} f(x)_{k_2}$ **then**
11         return True
12      **else**
13         return False

---

image. Different from [33, 34] that focus on *semantic* distortion, our method addresses *visual* distortion and aims to generate adversarial examples that are visually indistinguishable from the original image.

## 3   METHOD

### 3.1   *Learning Noise Distribution Based on Visual Distortion*

An attack model is an adversary that constructs adversarial examples against certain networks. Let $f : x \rightarrow f(x)$ be the target network that accepts an input $x \in \mathbb{R}^n$ and produces an output $f(x) \in \mathbb{R}^m$. $f(x)$ is a vector and $f(x)_k$ represents its $k$th entry, denoting the score of the $k$th class. $y = argmax_k f(x)_k$ is the predicted class. Given a valid input $x$ and the corresponding predicted class $y$, an adversarial example [35] $x'$ is similar to $x$ yet results in an incorrect prediction $argmax_k f(x')_k \neq y$. In an additive attack, an adversarial example $x'$ is a perturbed input with additive noise $\delta$ such that $x' = x + \delta$. The problem of generating an adversarial example is equivalent to produce noise map $\delta$ that causes wrong prediction for the perturbed input. Thus a successful attack is to find $\delta$ such that $argmax_k f(x + \delta)_k \neq y$. Since this constraint is highly non-linear, the loss function is usually rephrased in a different form [5]:

$$L(x, x + \delta) = \max(0, f(x + \delta)_y - \max_{k \neq y} f(x + \delta)_k) \quad (1)$$

The attack is successful when $L = 0$. It's noted that such a loss does not take the visual impact into consideration, for which the adversarial example could suffer from significant visual distortion. In order to constrain the visual distortion caused by the difference between $x$ and $x + \delta$, we adopt a perceptual distance metric $d(x, x+\delta)$ into the loss function with a predefined hyperparameter $\lambda$:

$$
\begin{aligned}
L(x, x + \delta) = &\max(0, f(x + \delta)_y - \max_{k \neq y} f(x + \delta)_k) \\
&+ \lambda d(x, x + \delta)
\end{aligned}
\quad (2)
$$

where smaller $d(x, x + \delta)$ indicates less visual distortion. $d$ can be any form of metric that measures the perceptual distance between $x$ and $x + \delta$, such as well-established $1 - \text{SSIM}$ [36] or

LPIPS [37]. $\lambda$ manages the trade-off between a successful attack and the visual distortion caused by the attack. The effects of $\lambda$ will be further discussed in Section 4.1.

Minimizing the above loss function faces a challenge that $L$ is not differentiable since the black-box adversary does not have access to the gradients of $L$ and the metric $d(x, x + \delta)$ might be calculated in a non-differentiable way. To address this problem, we *explicitly* assume a flexible noise distribution of $\delta$ in the discrete space, in the sense that the noise values and their probabilities are discrete. And the gradient of $L$ can be estimated by sampling from this distribution. Suppose that $\delta$ follows a distribution $p_\theta$ parameterized by $\theta$, *i.e.*, $\delta \sim p_\theta$. For the $j$th pixel in an image, we make its noise distribution $p_{\theta^j} = \text{softmax}(\theta^j)$, where $\theta^j$ is a vector and each element in it denotes a probability value. By sampling noise from the distribution $p_\theta$, $\theta$ can be learned to minimize the expectation of the above loss such that the attack is successful (*i.e.*, alters the predicted label) and the produced adversarial example is less distorted (*i.e.*, small $d$):

$$\text{minimize } \mathbb{E}_{\delta \sim p_\theta}[L(x, x + \delta)] \quad (3)$$

For the $j$th pixel, we define its noise's sample space to be a set of discrete values ranging from $-\epsilon$ to $\epsilon$: $\delta^j \in \{\epsilon, \epsilon - \frac{\epsilon}{N}, \epsilon - 2\frac{\epsilon}{N}, ..., 0, ... - \epsilon\}$, where $N$ is the sampling frequency and $\frac{\epsilon}{N}$ is the sampling interval. The noise value $\delta^j$ of the $j$th pixel is sampled from this sample space by following $p_{\theta^j}$, $p_{\theta^j} \in \mathbb{R}^{2N+1}$.

Given $W$ and $H$ the width and height of an image, respectively, since each pixel has its own noise distribution $p_{\theta^j}$ of length $2N + 1$, the number of parameters for the entire image is $(2N + 1)WH$. Note that we do not consider the difference of color channels in order to reduce the size of the sample space. Otherwise the number of parameters would be tripled. Thus, the same noise value is sampled for each RGB channel of a pixel. To estimate $\theta$, we adopt policy gradient [38] to make the above expectation differentiable with respect to $\theta$. Using REINFORCE, we have the differentiable loss function $F(\theta)$:

$$
\begin{aligned}
F(\theta) &= \mathbb{E}_{\delta \sim p_\theta}[L(x, x + \delta) - b] \\
&= (L(x, x + \delta) - b) \log(p_\theta(\delta))
\end{aligned}
\quad (4)
$$

$$
\begin{aligned}
\nabla F(\theta) &= \nabla_\theta \mathbb{E}_{\delta \sim p_\theta}[L(x, x + \delta) - b] \\
&= (L(x, x + \delta) - b)(1 - p_\theta(\delta))
\end{aligned}
\quad (5)
$$

where $b$ is introduced as a *baseline* in the expectation with specific meaning: 1) when $L(x, x + \delta) < b$, the sampled noise map $\delta$ returns low $L$, and its probability $p_\theta(\delta)$ increases through gradient descent; 2) when $L(x, x + \delta) = b$, $\nabla F(\theta) = 0$ and $p_\theta(\delta)$ remains unchanged; 3) when $L(x, x + \delta) > b$, the sampled noise map $\delta$ returns high $L$, and its probability $p_\theta(\delta)$ decreases through gradient descent. To sum up, $L(x, x + \delta)$ is forced to improve over $b$. At the iteration $t$, we choose $b = \min_{T=0,1,...t-1} L(x, x+\delta_T)$ such that $L$ improves over the obtained minimal loss.

The above expectation is estimated using a single Monte Carlo sampling at each iteration, and the sampling of noise map $\delta$ is critical. Simply sampling $\delta_t$ at the iteration $t$ on the entire image might cause large variance on the norm of the noise, *i.e.*, $\|\delta_t - \delta_{t-1}\|_2$. Therefore, to ensure a small variance, with $T^* = \text{argmin}_{T=0,1,...t-1} L(x, x+\delta_T)$, only a number of $qWH$ pixels' noise values are resampled in $\delta_{T^*}$, while $(1-q)WH$ pixels' noise values remain unchanged:

$$\delta_{t+1} \leftarrow \text{resample}(\delta_{T^*}, q; \delta_t)_{\delta_t \sim p_{\theta_t}} \quad (6)$$
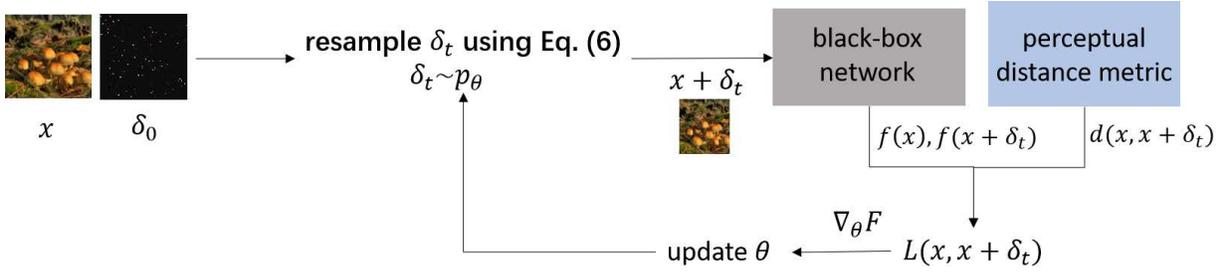
3

Figure 2: Framework of the proposed attack.

The above equation means replacing $qWH$ pixels' noise values in noise map $\delta_{T^*}$ with those in $\delta_t$, which are sampled from distribution $p_{\theta_t}$. In other words, if $q = 0.01$, then only a random 1% of $\delta_{T^*}$ is updated at each iteration. As shown in Fig. 2, after sampling $\delta_t$, the feedback $L(x, x + \delta_t)$ from the black-box and the perceptual distance metric decide the update of the distribution $p_{\theta_t}$. The iteration stops when the attack is successful, i.e., $\max(0, f(x + \delta_t)_y - \max_{k \neq y} f(x + \delta_t)_k) = 0$.

### 3.2 Proof of Convergence

Ruan *et al.* [39] shows that feed-forward DNNs (Deep Neural Networks) are Lipschitz continuous with a Lipschitz constant $K$, for which we have

$$\forall t, \|f(x + \delta_t) - f(x + \delta_{T^*})\|_2 \leq K\|\delta_t - \delta_{T^*}\|_2 \quad (7)$$

At the iteration $t$, since only a small part of the noise map is updated, it can be assumed that

$$|\max_{k \neq y} f(x + \delta_t)_k - \max_{k \neq y} f(x + \delta_{T^*})_k| \leq C \quad (8)$$

where $C$ is a constant. Suppose that the perceptual distance metric $d$ is normalized to $[0, 1]$. Substituting the inequalities (7) and (8) in our definition of $L$ in Eq. (2) gets the following inequalities:

$$\begin{aligned} &|L(x, x + \delta_t) - L(x, x + \delta_{T^*})| \\ &\leq K\|\delta_t - \delta_{T^*}\|_2 + C + \lambda \\ &\leq 2KWH\epsilon cq + C + \lambda \end{aligned} \quad (9)$$

Ideally, $L(x, x + \delta_t) - L(x, x + \delta_{T^*})$ accurately quantifies the difference of the perturbed image even when only one noise value for just a single pixel at the iteration $t$ is different from that at $T^*$. Let $\delta^{ij}$ represent a special noise map, whose $j_{th}$ pixel's noise value is the $i_{th}$ element in its sample space and the other pixels' noise values are 0. Note that the length of the sample space for each pixel is $2N + 1$. Similarly, $p_{\theta_t}(\delta^{ij})$ denotes the probability of the $i_{th}$ element in the sample space of the $j_{th}$ pixel. By sampling every element in the sample space of the $j_{th}$ pixel, we define $l_t^j$ and $p_{\theta_t^j}$ to be a vector:

$$\forall j \in \{1, 2, ..., WH\}, l_t^j = vector[L(x, x + \delta^{ij}) - L(x, x + \delta_{T^*})], \\ i = 1, 2, ..., 2N + 1 \quad (10)$$

$$\forall j \in \{1, 2, ..., WH\}, p_{\theta_t^j} = vector[p_{\theta_t}(\delta^{ij})], \\ i = 1, 2, ..., 2N + 1 \quad (11)$$

Although the above equations are only meaningful under the ideal situation where $L$ can quantify the difference of just one

perturbed pixel, we use these equations for a theoretical proof of convergence. In the ideal situation, the gradient of the $j_{th}$ pixel's parameters can be calculated exactly as

$$\nabla F(\theta_t^j) = l_t^j \cdot (\mathbf{1} - p_{\theta_t^j}) \quad (12)$$

According to Eq. (9) when the number of the resampled pixels $qWH = 1$, we have

$$|L(x, x + \delta^{ij}) - L(x, x + \delta_{T^*})| \leq 2K\epsilon c + C + \lambda \quad (13)$$

Note that for $\forall t_1, t_2$ that share the same $T^*$, $l_{t_1}^j$ is equal to $l_{t_2}^j$. Thus, using Eq. (13), we have

$$\begin{aligned} &\|\nabla F(\theta_{t_1}^j) - \nabla F(\theta_{t_2}^j)\|_2 \\ &\leq (2N + 1)(2K\epsilon c + C + \lambda)\|\text{softmax}(\theta_{t_1}^j) - \text{softmax}(\theta_{t_2}^j)\|_2 \end{aligned} \quad (14)$$

In practice, we adopt a single Monte Carlo sampling instead of sampling every noise values for every pixel, for which $2N + 1$ should be replaced by 1 in the above inequality. The inequality (14) thus becomes:

$$\begin{aligned} &\|\nabla F(\theta_{t_1}^j) - \nabla F(\theta_{t_2}^j)\|_2 \\ &\leq (2K\epsilon c + C + \lambda)\|\text{softmax}(\theta_{t_1}^j) - \text{softmax}(\theta_{t_2}^j)\|_2 \\ &\leq (2K\epsilon c + C + \lambda)\|\theta_{t_1}^j - \theta_{t_2}^j\|_2 \end{aligned} \quad (15)$$

The standard softmax function disappears because it is Lipschitz continuous with the Lipschitz constant being 1 [40]. Finally, we have the inequality for $\|\nabla F(\theta_{t_1}) - \nabla F(\theta_{t_2})\|_2$:

$$\|\nabla F(\theta_{t_1}) - \nabla F(\theta_{t_2})\|_2 \leq (2K\epsilon c + C + \lambda)\|\theta_{t_1} - \theta_{t_2}\|_2 \quad (16)$$

The above inequality proves that $F(\theta)$ is $L$-smooth with the Lipschitz constant being $2K\epsilon c + C + \lambda$. If $F(\theta)$ is convex, the exact number of steps that Stochastic Gradient Descent (SGD) takes to convergence is $\frac{(2K\epsilon c + C + \lambda) \cdot \|\theta_0 - \theta^*\|_2^2}{\xi}$, where $\xi$ is an arbitrary small tolerable error ($\xi > 0$). However, since deep network $L$ is usually highly non-convex, we need to consider the situation where $F(\theta)$ is non-convex.

Let the SGD update be

$$\theta_{t+1} = \theta_t + \eta_t g(\theta_t) \quad (17)$$

where $\eta_t$ is the learning rate and $g(\theta_t)$ is the stochastic gradient. We assume that the variance of the stochastic gradient is upper bounded by $\sigma^2$:

$$\mathbb{E}[\|\nabla F(\theta) - g(\theta)\|_2^2] \leq \sigma^2 < \infty \quad (18)$$

4

Table 1: Ablation results of the perceptual distance metric, $\lambda$ and sampling frequency $N$. Smaller $1-$SSIM, LPIPS and CIEDE2000 indicate less visual distortion.

| Sampling Frequency | Perceptual Metric | $\lambda$ | Success Rate | $1 -$ SSIM | LPIPS | CIEDE2000 | Avg. Queries |
|---|---|---|---|---|---|---|---|
| | - | 0 | **100%** | 0.091 | 0.099 | 0.941 | **356** |
| $N = 1$ | $1 -$ SSIM | 10 | **100%** | 0.076 | 0.081 | 0.741 | 401 |
| | | 100 | 97.4% | 0.036 | 0.051 | 0.703 | 1395 |
| | | 200 | 92.2% | 0.025 | 0.040 | 0.622 | 2534 |
| | | dynamic | **100%** | **0.009** | 0.009 | **0.204** | 7678 |
| | LPIPS | 10 | **100%** | 0.080 | 0.078 | 0.762 | 450 |
| | | 100 | 98.1% | 0.049 | 0.052 | 0.711 | 1174 |
| | | 200 | 95.1% | 0.038 | 0.045 | 0.635 | 1928 |
| | | dynamic | **100%** | 0.015 | **0.005** | 0.277 | 6694 |
| None | $1 -$ SSIM | 10 | **100%** | 0.118 | 0.142 | 5.936 | 426 |
| $N = 2$ | $1 -$ SSIM | 10 | 99.7% | 0.071 | 0.074 | 0.846 | 520 |
| $N = 5$ | $1 -$ SSIM | 10 | 99.5% | 0.069 | 0.070 | 0.877 | 665 |
| $N = 10$ | $1 -$ SSIM | 10 | 98.7% | 0.062 | 0.075 | 0.879 | 669 |
| $N = 12$ | $1 -$ SSIM | 10 | 98.7% | 0.071 | 0.075 | 0.882 | 673 |

And we select $\eta_t$ that satisfies

$$\sum_{t=1}^{\infty} \eta_t = \infty \text{ and } \sum_{t=1}^{\infty} \eta_t^2 < \infty \quad (19)$$

Condition (19) can be easily satisfied with a decaying learning rate, e.g., $\eta_t = \frac{1}{\sqrt{t}\ln(t+1)}$. According to Lemma 1 and Theorem 2 in [41], using the $L$-smooth property of $F(\theta)$, $\nabla F(\theta_t)$ goes to 0 with probability 1. This means that with probability 1 for any $\xi > 0$ there exists $N_\xi$ such that $\nabla F(\theta_t) \leq \xi$ for $t \geq N_\xi$. Unfortunately, unlike in the convex case, we do not know the exact number of steps that SGD takes to convergence.

The above proof simply aims to theoretically show that the proposed method converges in finite steps, although possibly in a rather slow speed. From the "Avg. Queries" in the following experiments, we can see that the actual computational cost is affordable and comparable to some of the query-efficient attacks.

## 4 EXPERIMENTS

Following previous work [25, 8], we validate the effectiveness of our model on the large-scale ImageNet [42] dataset. We use three pretrained classification networks on Pytorch as the black-box networks: InceptionV3 [43], ResNet50 [44] and VGG16bn [45]. The attack is performed on images that were correctly classified by the pretrained network. We randomly select 1000 images in the validation set for test, and all images are normalized to [0, 1]. We quantify our success in terms of the perceptual distance ($1-$SSIM, LPIPS and CIEDE2000) as we address the visual distortion caused by the attack. In these metrics, $1-$SSIM [36] measures the degradation of structural information in the adversarial examples. LPIPS [37] evaluates the perceptual similarity of two images with their normalized distance between their deep features. CIEDE2000 [46] measures perceptual color distance, which is developed by the CIE (International Commission on Illumination). Smaller value of these metrics denotes less

visual distortion. Except for $1-$SSIM, LPIPS and CIEDE2000, the success rate and average number of queries are also reported as in most previous work. The average number of queries refers to the average number of requests to the output of the black-box network.

We initialize the noise distribution $p_\theta$ to be a uniform distribution and noise $\delta_0$ to be 0. The learning rate is 0.01 and $q$ is set to be 0.01. In addition, we specify the shape of the resampled noise at each iteration to be a square [25, 24, 26], and adopt the tiling trick [8, 25] with tile size= 2. The upper bound $\epsilon$ of our attack is set to be 0.05 as in previous work.

### 4.1 Ablation Studies

In the ablation studies, the maximum number of queries is set to be 10, 000. The results are averaged on 1000 test images. In the following, we discuss the trade-off between visual distortion and query efficiency, the effects of using different perceptual distance metrics in the loss function, the results on different sampling frequencies and the influence of predefining a specific form of noise distribution.

**Trade-off between visual distortion and query efficiency.** Under the same $l_\infty$ ball, a query-efficient way to produce an adversarial example is to perturb most pixels with the maximum noise values $\pm\epsilon$ [24, 26]. However, such attack introduces large visual distortion, which could make the distorted image very annoying. To constrain the visual distortion, the perturbed pixels should be those who cause smaller visual difference while performing a valid attack, which takes extra queries to find. This brings the trade-off between visual distortion and query efficiency, which can be controlled by $\lambda$ in our loss function. As shown in Table 1, when $N = 1$ and $\lambda = 0$, the adversary does not consider visual distortion at all, and perturbs each pixel that is helpful for misclassification until the attack is successful. Thus, it causes the largest perceptual distance (0.091, 0.099 and 0.941)

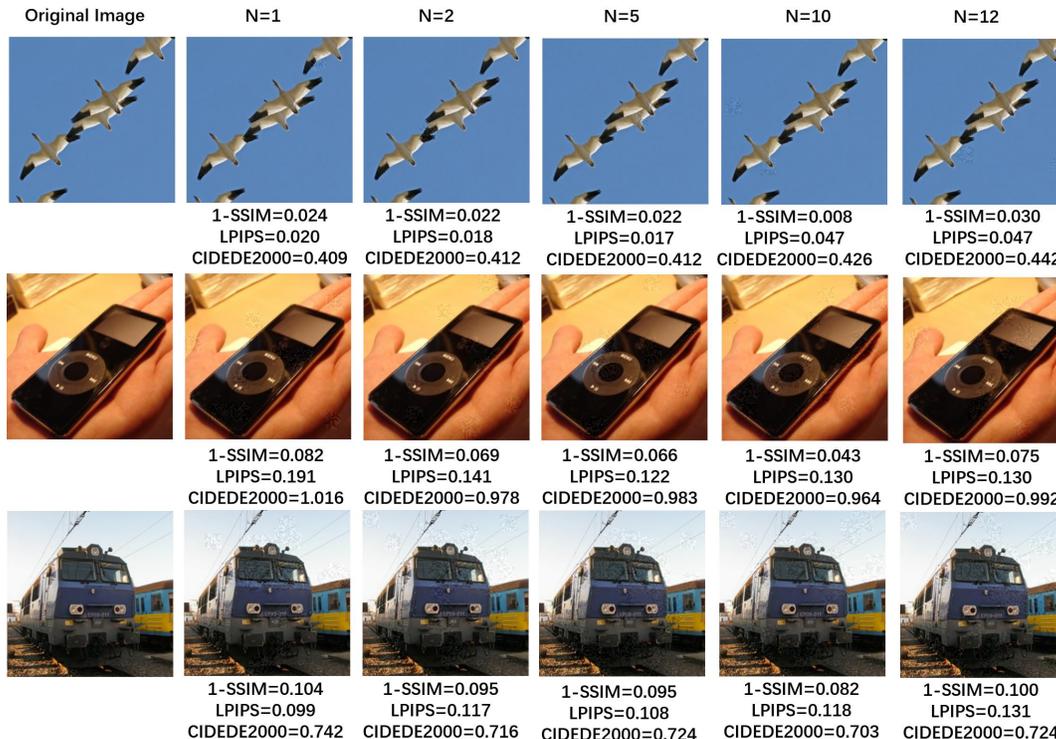| Original Image | N=1 | N=2 | N=5 | N=10 | N=12 |
|---|---|---|---|---|---|
| | 1-SSIM=0.024 LPIPS=0.020 CIDEDE2000=0.409 | 1-SSIM=0.022 LPIPS=0.018 CIDEDE2000=0.412 | 1-SSIM=0.022 LPIPS=0.017 CIDEDE2000=0.412 | 1-SSIM=0.008 LPIPS=0.047 CIDEDE2000=0.426 | 1-SSIM=0.030 LPIPS=0.047 CIDEDE2000=0.442 |
| | 1-SSIM=0.082 LPIPS=0.191 CIDEDE2000=1.016 | 1-SSIM=0.069 LPIPS=0.141 CIDEDE2000=0.978 | 1-SSIM=0.066 LPIPS=0.122 CIDEDE2000=0.983 | 1-SSIM=0.043 LPIPS=0.130 CIDEDE2000=0.964 | 1-SSIM=0.075 LPIPS=0.130 CIDEDE2000=0.992 |
| | 1-SSIM=0.104 LPIPS=0.099 CIDEDE2000=0.742 | 1-SSIM=0.095 LPIPS=0.117 CIDEDE2000=0.716 | 1-SSIM=0.095 LPIPS=0.108 CIDEDE2000=0.724 | 1-SSIM=0.082 LPIPS=0.118 CIDEDE2000=0.703 | 1-SSIM=0.100 LPIPS=0.131 CIDEDE2000=0.724 |

Figure 3: Adversarial examples under different sampling frequency. From left to right is the original image, the adversarial examples from $N = 1, 2, 5, 10, 12$, respectively.

with the least number of queries (356). As $\lambda$ increases to 200, all the perceptual metrics decrease at the cost of more queries and lower success rate. The maximum $\lambda$ in Table 1 is 200 since further increasing it causes the success rate to be lower than 90%. In addition, as in [17], we perform a dynamic line search on the choice of $\lambda$ to see the best perceptual scores the adversary can achieve, where $\lambda \in [0, 1000]$. Comparing with fixed $\lambda$ values, using dynamic values of $\lambda$ greatly boosts the performance on perceptual metrics with 100% attack success rate, at the cost of dozens of times the number of queries. Fig. 4 gives several visualized examples on different $\lambda$, where adversarial examples with larger $\lambda$ suffer from less visual distortion.

**Ablation studies on the perceptual distance metric.** The perceptual distance metric $d$ in the loss function is predefined to measure the visual distortion between the adversarial example and the original image. We adopt $1 - SSIM$ and LPIPS as the perceptual distance metric to optimize, respectively, and report their results in Table 1. When $\lambda = 10$, optimizing $1 - SSIM$ shows better score on $1 - SSIM$ (0.076 v.s. 0.080) and CIEDE2000 (0.721 v.s. 0.742) whilst optimizing LPIPS has better performance on LPIPS (0.078 v.s. 0.081). However, when $\lambda$ increases to 100 and 200, optimizing $1 - SSIM$ gives better scores on both $1 - SSIM$ and LPIPS. Therefore, we set the perceptual distance metric to be $1 - SSIM$ in the following experiments.

**Sampling frequency.** Sampling frequency decides the size of the sample space of $\delta$. Setting higher frequency means there are more noise values to explore through sampling. In Table 1, increasing the sampling frequency from $N = 1$ to $N = 2$

reduces the perceptual distance to some extent at the cost of lower success rate. On the other hand, further increasing $N$ to 12 does not essentially reduce the distortion yet lowers the success rate. We set the sampling frequency $N = 1$ in the following experiments. Note that the maximum sampling frequency is $N = 12$ because the sampling interval in RGB color space (*i.e.*, $255 * 0.05/N$) would be less than 1 if $N > 12$. See Fig. 3 for a few adversarial examples from different sampling frequencies.

**Noise Distribution.** In the proposed algorithm, we adopt a flexible noise distribution instead of predefining it to be a specific form. Therefore, we conducted ablation studies on assuming the distribution to be a regular form as in NAttack [21]. Specifically, we let the noise distribution be an isometric normal distribution while $\lambda = 10$ in the loss function, and perform attacks by estimating mean and variance as Eq. (10) in [21]. As reported in the tenth row of Table 1, under the same experimental setting, it is clear that fixing the noise distribution to be a specific isometric normal distribution degrades the overall performance. We think it is because the distribution that minimizes the perceptual distance is unknown, which might not follow a Guassian distribution or other regular form of distribution. To approximate an unknown distribution, it is better to allow the noise distribution to be a free form as in the proposed approach, and let it be learned by minimizing the perceptual distance.

### 4.2 Out-of-Object Attack

Most existing classification networks [44, 47] are based on Convolutional Neural Network (CNN), which gradually aggregates
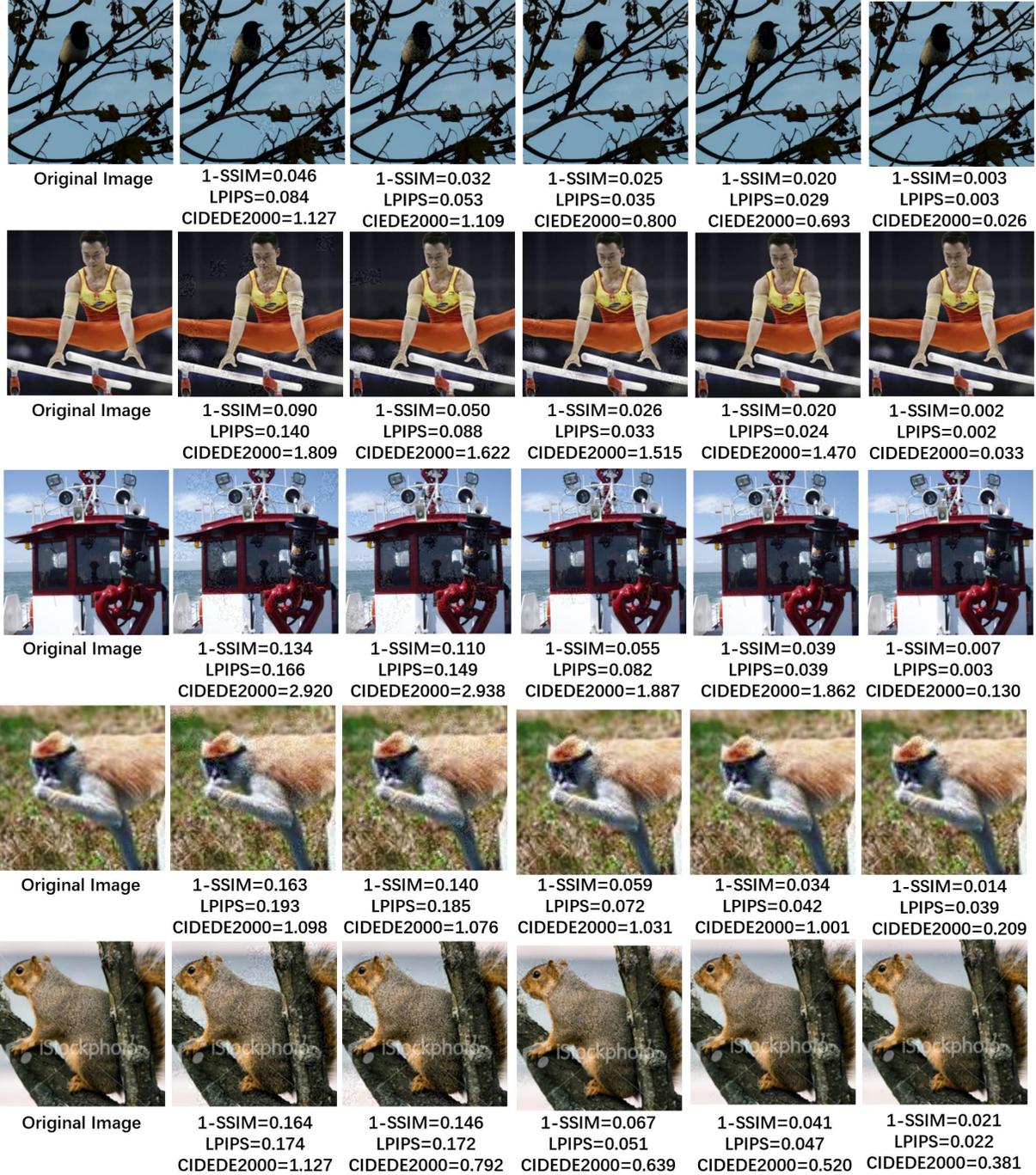
Figure 4: Visualized examples of the proposed attack. From left to right is the original image, the adversarial examples on $\lambda = 0, \lambda = 10, \lambda = 100, \lambda = 200$, dynamic $\lambda$, respectively.

Table 2: Results of the out-of-object attack on ImageNet when $\lambda = 10, N = 1$ and the perceptual distance metric being $1 - $ SSIM. I, R and V represent InceptionV3, ResNet50 and VGG16bn, respectively.

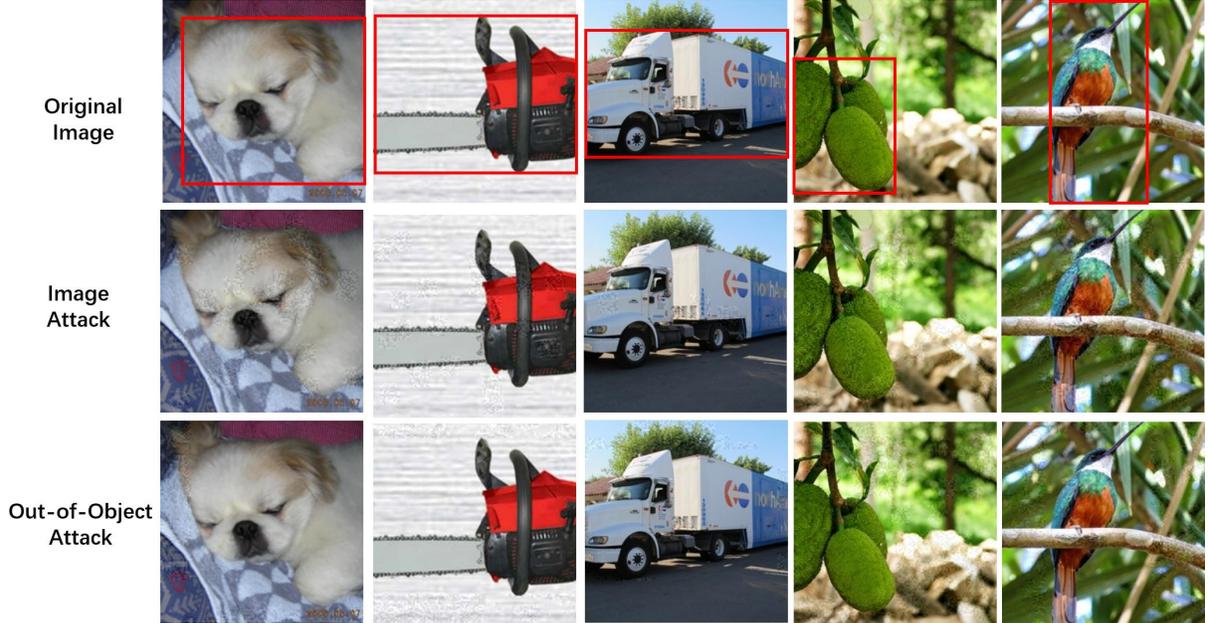| Attacked Range | Success Rate | | | $1 - $ SSIM | | | LPIPS | | | CIEDE2000 | | | Avg. Queries | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | I | R | V | I | R | V | I | R | V | I | R | V | I | R | V |
| Image | **100%** | **100%** | **100%** | 0.078 | 0.076 | **0.072** | 0.096 | 0.081 | 0.079 | 0.692 | **0.741** | 0.699 | **845** | **401** | **251** |
| Out-of-object | 90.1% | 93.8% | 94.7% | **0.071** | **0.069** | 0.074 | **0.081** | **0.065** | **0.070** | **0.678** | 0.805 | **0.687** | 4275 | 3775 | 3104 |

Figure 5: Visualized adversarial examples in out-of-object attack. The red bounding box locates the target object in the original image. In *out-of-object attack*, the adversary is only allowed to perturb pixels that are out of the object bounding box. In *image attack*, the adversary can perturb any pixel in the image.

contextual information in deeper layers. Therefore, it is possible to fool the classifier by just attacking the "context", *i.e.*, background that is out of the target object. Attacking just the out-of-object pixels constrains the number and the position of pixels that can be perturbed, which might further reduce the visual distortion caused by the noise. To locate the object in a given image, we exploited the object bounding box provided by ImageNet. An out-of-object mask is then created according to the bounding box such that the model is only allowed to attack pixels that are out of the object, as shown in Fig. 5. In Table 2, we report results of InceptionV3, ResNet50 and VGG16bn with the maximum queries= 40,000. The attack is performed on images whose masks are at least 10% large of the image area. The results show that attacking just the out-of-object pixels can also cause misclassification of the object with over 90% success rate. Compared with image attack, the out-of-object attack is more difficult for the adversary in that it requires more number of queries (4275/3775/3104) yet has lower success rate (90.1%/93.8%/94.7%). On the other hand, the out-of-object attack indeed reduces visual distortion of the adversarial examples on the three networks.

Table 3: Comparison of the undefended (v3) and defended (v3$_{adv-ens4}$) InceptionV3. The defended InceptionV3 adopts ensemble adversarial training.

| Network | Clean Accuracy | After Attack | 1 − SSIM | LPIPS | CIEDE2000 | Avg. Queries |
|---------|----------------|--------------|----------|-------|-----------|--------------|
| v3 | 75.8% | 0.8% | 0.096 | 0.149 | 0.862 | 531 |
| v3$_{adv-ens4}$ | 73.4% | 1.8% | 0.103 | 0.154 | 0.979 | 771 |

### 4.3 Attack Effectiveness on Defended Network

In the above experiments, we show that our black-box model can attack the *undefended* network with high success rate. To

evaluate the strength of the proposed attack in *defended* situation, we further attack the InceptionV3 network that adopts ensemble adversarial training (*i.e.*, v3$_{adv-ens4}$). Following [48], we set $\epsilon = 0.0625$ and randomly select 10,000 images from the ImageNet validation set for test. The maximum number of queries is 10,000. The performance of the attacked network is reported in Table 3, where clean accuracy is the classification accuracy before attack. Note that v3 is slightly different from InceptionV3 in Table 1 in that the pretrained model of v3 comes from Tensorflow, which is the same platform of the pretrained model of v3$_{adv-ens4}$. Compared with undefended network, attacking defended one causes larger visual distortion. However, the proposed attack can still reduce the classification accuracy from 73.4% to 1.8%, which demonstrates its effectiveness under defend.

### 4.4 Comparison with Other Attacks

Since our approach addresses improving the visual similarity between the adversarial example and the original image, it might cost more number of queries to construct a less distorted adversarial example. To show that such costs are affordable, we compare our attack to recently proposed black-box attacks: Sign-Hunter [9], NAttack [21], AutoZOOM [17], Bandits [8], Square Attack [26] and TREMBA [20]. For fair comparison, in Table 4, methods marked with -SSIM and **Ours** introduce $\lambda \cdot (1 - SSIM)$ to the loss function with $\lambda = 10$. Note that AutoZOOM performs line search on the choice of $\lambda$, for which we adopt the same strategy and denotes this variant of our method as **Ours**($\lambda_{dynamic}$). The results of the above methods are reproduced using the official codes provided by the authors. We use the default parameter settings of the corresponding attack, and set the maximum number of queries to be 10,000. See Table 5 for the experimental settings of different methods. In Table 4, Comparing ap-
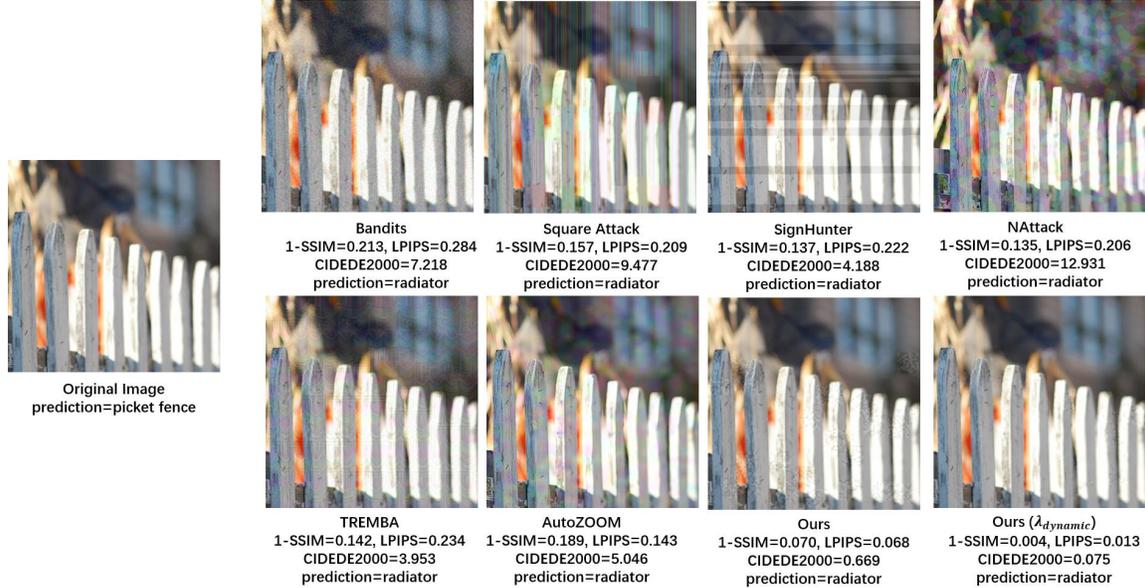
8

Figure 6: Adversarial examples from different attacks with perceptual distance scores.

Table 4: Results of different attacks on ImageNet. I, R and V represent InceptionV3, ResNet50 and VGG16bn, respectively.

| Attack | Success Rate | | | $1 - SSIM$ | | | LPIPS | | | CIEDE2000 | | | Avg. Queries | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | I | R | V | I | R | V | I | R | V | I | R | V | I | R | V |
| SignHunter [9] | 98.4% | - | - | 0.157 | - | - | 0.117 | - | - | 3.837 | - | - | 450 | - | - |
| NAttack [21] | 99.5% | - | - | 0.133 | - | - | 0.212 | - | - | 5.478 | - | - | 524 | - | |
| AutoZOOM [17] | 100% | - | - | 0.038 | - | - | 0.059 | - | - | 3.33 | - | - | 1010 | - | - |
| Bandits [8] | 96.5% | 98.8% | 98.2% | 0.343 | 0.307 | 0.282 | 0.201 | 0.157 | 0.140 | 8.383 | 8.552 | 8.194 | 935 | 705 | 388 |
| Square Attack [26] | 99.7% | **100%** | **100%** | 0.280 | 0.279 | 0.299 | 0.265 | 0.243 | 0.247 | 9.329 | 9.425 | 9.429 | **237** | **62** | **30** |
| TREMBA [20] | 99.0% | **100%** | 99.8% | 0.161 | 0.161 | 0.160 | 0.188 | 0.189 | 0.187 | 4.413 | 4.400 | 4.421 | - | - | - |
| SignHunter-SSIM | 97.6% | - | - | 0.220 | - | - | 0.157 | - | - | 3.832 | - | - | 642 | - | - |
| NAttack-SSIM | 97.3% | - | - | 0.128 | - | - | 0.210 | - | - | 5.021 | - | - | 666 | - | - |
| AutoZOOM-SSIM | **100%** | - | - | 0.028 | - | - | 0.048 | - | - | 2.98 | - | - | 2245 | - | - |
| Bandits-SSIM | 80.0% | 89.3% | 89.7% | 0.333 | 0.303 | 0.275 | 0.200 | 0.163 | 0.135 | 8.838 | 8.666 | 8.194 | 1318 | 1020 | 793 |
| Square Attack-SSIM | 99.2% | 100% | 100% | 0.260 | 0.268 | 0.292 | 0.256 | 0.238 | 0.245 | 9.301 | 9.462 | 9.451 | 278 | 65 | **30** |
| TREMBA-SSIM | 98.5% | 100% | 99.8% | 0.160 | 0.160 | 0.159 | 0.185 | 0.186 | 0.183 | 4.410 | 4.396 | 4.421 | - | - | - |
| Ours | 98.7% | **100%** | **100%** | 0.075 | 0.076 | 0.072 | 0.094 | 0.081 | 0.079 | 0.692 | 0.741 | 0.699 | 731 | 401 | 251 |
| Ours($\lambda_{dynamic}$) | **100%** | **100%** | **100%** | **0.016** | **0.009** | **0.006** | **0.023** | **0.009** | **0.005** | **0.215** | **0.204** | **0.155** | 7311 | 7678 | 7620 |

proaches that use fixed $\lambda$ value (i.e., Signhunter-SSIM, NAttack-SSIM, Bandits-SSIM, Square Attack-SSIM, TREMBA-SSIM, AdvGAN-SSIM and Ours), we can see that the proposed method outperforms other attacks on reducing perceptual distance, while the average number of queries is comparable to Bandits. On the other hand, Ours($\lambda_{dynamic}$) achieves state-of-the-art performance on 1-SSIM, LPIPS and CIEDE2000 when compared with methods that perform line search over $\lambda$ (i.e., AutoZOOM and AutoZOOM-SSIM). In general, except for Signhunter, introducing perceptual distance metric in the objective function helps reduce visual distortion in other attacks. The visualized adversarial examples from different attacks are given in Fig. 6, which shows that our model produces less distorted adversarial examples. More examples can be found in Fig. 7.

We noticed that adversarial examples from SignHunter have horizontal-stripped noise and Square Attack generates adversarial examples with vertical-stripped noise. Stripped noise is helpful in improving query efficiency since the classification

network is quite sensitive to such noise [26]. However, from the perspective of visual distortion, such noise greatly degrades the image quality. The adversarial examples of Bandits are relatively perceptible-friendly, but the perturbation affects most pixels in the image, which causes visually "noisy" effects, especially in a monocolor background. The noise maps from NAttack and AutoZOOM appear to be regular color patches all over the image due to their large tiling size in the methods.

We also conducted subjective studies for further validation. Specifically, we randomly chose two adversarial examples, where one is generated by our approach (Ours($\lambda_{dynamic}$)) and the other is given by the attacks (excluding ours) in Table 4. We show each human evaluator the two adversarial examples, and ask him/her which one is less distorted compared with the original image. Figure 8 gives an picture that we show to the evaluator. Note that the order of the two adversarial examples in the triplet is randomly permuted. We asked 10 human evaluators in total, each made judgements over 100 triplets of images. As a
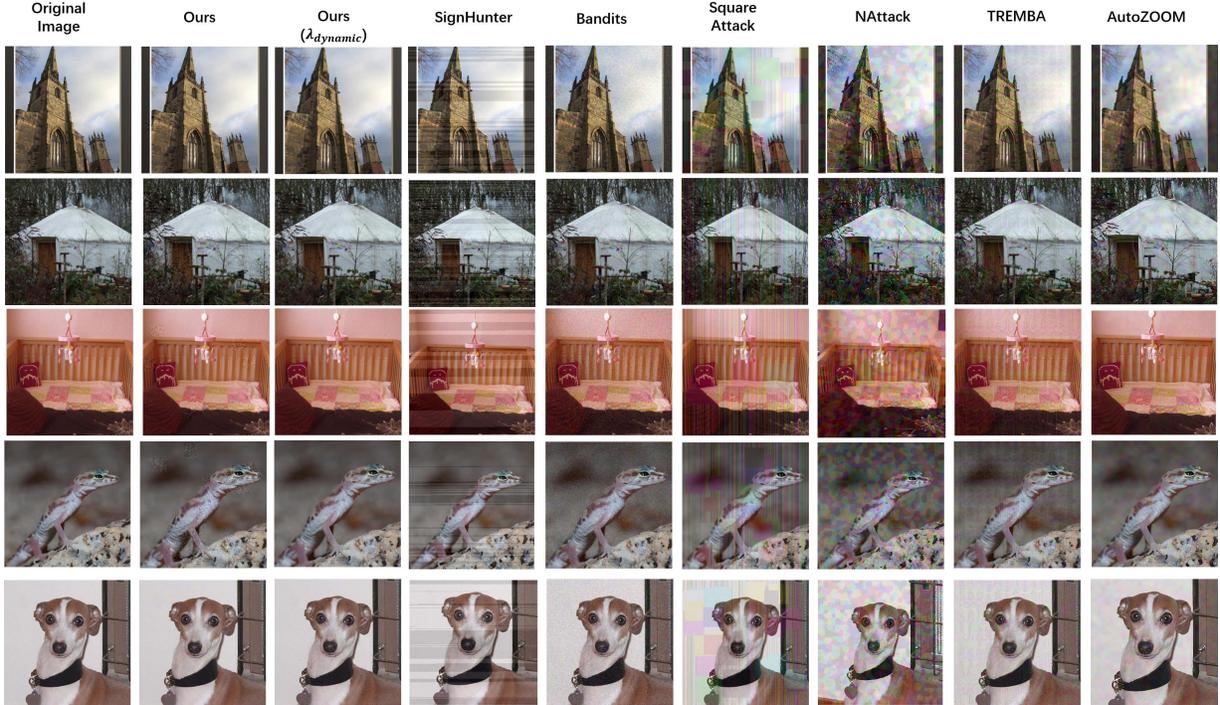
Figure 7: More visualized adversarial examples from different attacks.



Figure 8: An example of the pictures that we show to the evaluator. One of (a)(b) is produced by our model and the other is from the attacks (excluding ours) in Table 4.

Table 5: Experimental settings.

| Method | $\lambda$ | Max. Iterations |
|---|---|---|
| Signhunter-SSIM | 10 | 10,000 |
| NAttack-SSIM | 10 | 10,000 |
| AutoZOOM-SSIM | dynamic, $\lambda \in [0, 1000]$ | 10,000 |
| Bandits-SSIM | 10 | 10,000 |
| Square Attack-SSIM | 10 | 10,000 |
| TREMBA-SSIM | 10 | - |
| Ours | 10 | 10,000 |
| Ours($\lambda_{dynamic}$) | dynamic, $\lambda \in [0, 1000]$ | 10,000 |

result, adversarial examples generated by our method are thought to have less noticeable noise 82.1% of the time, while 10.0% of the time the evaluators think both examples are distorted at the same level. Therefore, the subjective results further prove that the proposed method effectively reduces visual distortion in adversarial examples.

### 4.5 Other $l_p$ Attacks

Although our method in this paper is based on $l_\infty$ attack, the perceptual distance metric $d$ in the loss function can be replaced by other $l_p$ ($p = 0, 1, 2$) distance. We did not discuss it in the above experiments because these $l_p$ distance metrics are less accurate in measuring the *perceptual* distance between images compared to the specifically designed metrics, such as well-established $1 - $ SSIM and LPIPS. Nevertheless, we still present the results of other $l_p$ ($p = 0, 1, 2$) attacks in Table 6, where the $l_p$ distance is normalized to $[0, 1]$ in the loss function. Specifically, $d(x, x + \delta) = \frac{l_p(x, x+\delta)}{\max_\delta(l_p(x, x+\delta))}$, where $l_p(x, x + \delta)$ is the $l_p$ distance between the original image $x$ and the perturbed image $x+\delta$. As in the paper, we set $\lambda = 10, \epsilon = 0.05$ and the maximum number of queries being $10,000$. We find that the raw $l_0$ and $l_1$ scores have much higher order of magnitude compared with other metrics, and thus the normalized scores of $l_0$ and $l_1$ distances are reported in Table 6. Note that when the sampling frequency $N = 1$, $l_0$ distance is equivalent to $l_1$ distance in that

$$
\begin{aligned}
\frac{l_1(x, x + \delta)}{\max_\delta(l_1(x, x + \delta))} &= \frac{mc \cdot \epsilon}{WHc \cdot \epsilon} \\
&= \frac{m}{WH} \\
&= \frac{l_0(x, x + \delta)}{\max_\delta(l_0(x, x + \delta))}
\end{aligned}
\tag{20}
$$

where $m$ is the number of perturbed pixels. $W, H$ and $c$ are the width, height and number of channels of a given image, respectively. Table 6 shows that optimizing $l_0$ distance gives better performance on both the perceptual distance metrics and the $l_p$ distance metrics.

10

Table 6: Results of other $l_p$ attacks on ResNet50 when $\lambda = 10$. The raw $l_0$ and $l_1$ scores have much higher order of magnitude compared with other metrics, and thus the normalized scores of $l_0$ and $l_1$ distances are reported.

| Distance Metric | Sampling Frequency | Success Rate | $1 - \mathrm{SSIM}$ | LPIPS | CIEDE2000 | $l_0$ | $l_1$ | $l_2$ | Avg. Queries |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | **99.5%** | 0.077 | 0.083 | 0.795 | **0.133** | 0.130 | 6.75 | 536 |
| $l_0$ | 2 | 99.2% | 0.065 | 0.069 | **0.768** | 0.159 | 0.118 | 5.88 | 679 |
| | 5 | 97.9% | **0.058** | **0.065** | 0.789 | 0.177 | **0.118** | **5.19** | 960 |
| | 1 | **99.5%** | 0.077 | 0.083 | 0.795 | **0.133** | 0.130 | 6.75 | 536 |
| $l_1$ | 2 | **99.5%** | 0.070 | 0.076 | 0.773 | 0.176 | 0.130 | 6.14 | 658 |
| | 5 | 99.2% | 0.066 | 0.070 | 0.768 | 0.218 | 0.129 | 5.74 | 800 |
| | 1 | **99.5%** | 0.110 | 0.112 | 0.829 | 0.215 | 0.211 | 8.21 | **392** |
| $l_2$ | 2 | **99.5%** | 0.092 | 0.100 | 0.803 | 0.259 | 0.191 | 7.44 | 431 |
| | 5 | **99.5%** | 0.087 | 0.094 | 0.792 | 0.312 | 0.185 | 6.89 | 579 |

### 4.6 Conclusion

We introduce a novel black-box attack based on the induced visual distortion in the adversarial example. The quantified visual distortion, which measures the perceptual distance between the adversarial example and the original image, is introduced in our loss where the gradient of the corresponding non-differentiable loss function is approximated by sampling from a learned noise distribution. The proposed attack can achieve a trade-off between visual distortion and query efficiency by introducing the weighted perceptual distance metric in addition to the original loss. The experiments demonstrate the effectiveness of our attack on ImageNet as our model achieves much lower distortion when compared to existing attacks. In addition, it is shown that our attack is valid even when it's only allowed to perturb pixels that are out of the target object in a given image.

## REFERENCES

[1] H. Kwon, Y. Kim, H. Yoon, and D. Choi. Selective audio adversarial example in evasion attack on speech recognition system. *IEEE Transactions on Information Forensics and Security*, 15:526–538, 2020.

[2] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *Proc. International Conference on Learning Representations*, 2017.

[3] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In *Proc. ACM on Asia Conference on Computer and Communications Security*, 2017.

[4] N. Akhtar and A. Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access*, 6:14410–14430, 2018.

[5] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 2017.

[6] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *Pro. IEEE European Symposium on Security and Privacy, EuroS&P*, 2016.

[7] Matt Jordan, Naren Manoj, Surbhi Goel, and Alexandros G. Dimakis. Quantifying perceptual distortion of adversarial examples. *arXiv preprint arXiv:1902.08265*, 2019.

[8] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversarial attacks with bandits and priors. In *Proc. International Conference on Learning Representations*, 2019.

[9] Abdullah Al-Dujaili and Una-May O'Reilly. Sign bits are all you need for black-box attacks. In *Proc. International Conference on Learning Representations*, 2020.

[10] Y. Zhang, X. Tian, Y. Li, X. Wang, and D. Tao. Principal component adversarial example. *IEEE Transactions on Image Processing*, 29:4804–4815, 2020.

[11] Q. Zhang, K. Wang, W. Zhang, and J. Hu. Attacking black-box image classifiers with particle swarm optimization. *IEEE Access*, 7:158051–158063, 2019.

[12] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proc. International Conference on Learning Representations*, 2015.

[13] Cassidy Laidlaw and Soheil Feizi. Functional adversarial attacks. In *Proc. International Conference on Neural Information Processing Systems*, 2019.

[14] Tianhang Zheng, Changyou Chen, and Kui Ren. Distributionally adversarial attack. In *Proc. AAAI Conference on Artificial Intelligence*, 2019.

[15] Hongjun Wang, Guanbin Li, Xiaobai Liu, and Liang Lin. A hamiltonian monte carlo method for probabilistic adversarial attack and learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[16] Nicholas Carlini and David A. Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *ACM Workshop on Artificial Intelligence and Security*, 2017.

[17] Chun-Chen Tu, Pai-Shun Ting, Pin-Yu Chen, Sijia Liu, Huan Zhang, Jinfeng Yi, Cho-Jui Hsieh, and Shin-Ming Cheng. AutoZOOM: Autoencoder-based zeroth order optimization method for attacking black-box neural networks. In *Proc. AAAI Conference on Artificial Intelligence*, 2019.

[18] Yonggang Zhang, Ya Li, Tongliang Liu, and Xinmei Tian. Dual-path distillation: A unified framework to improve black-box attacks. In *Proc. International Conference on Machine Learning*, 2020.

[19] Huichen Li, Linyi Li, Xiaojun Xu, Xiaolu Zhang, Shuang Yang, and Bo Li. Nonlinear gradient estimation for query efficient blackbox attack.

[20] Zhichao Huang and Tong Zhang. Black-box adversarial attack with transferable model-based embedding. In *Proc. International Conference on Learning Representations*, 2019.

[21] Yandong Li, Lijun Li, Liqiang Wang, Tong Zhang, and Boqing Gong. NATTACK: learning the distributions of adversarial examples for an improved black-box attack on deep neural networks. In *Proc. International Conference on Machine Learning*, 2019.

[22] Shuyu Cheng, Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Improving black-box adversarial attacks with a transfer-based prior. In *Proc. International Conference on Neural Information Processing Systems*, 2019.

[23] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.

[24] Seungyong Moon, Gaon An, and Hyun Oh Song. Parsimonious black-box adversarial attacks via efficient combinatorial optimization. In *Proc. International Conference on Machine Learning*, 2019.

[25] Laurent Meunier, Jamal Atif, and Olivier Teytaud. Yet another but more efficient black-box adversarial attack: tiling and evolution strategies. *arXiv preprint arXiv:1910.02244*, 2019.

[26] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. *arXiv preprint arXiv:1912.00049*, 2019.

[27] Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. Generating adversarial examples with adversarial networks. 2018.

[28] Weijia Zhang. Generating adversarial examples in one shot with image-to-image translation gan. *IEEE Access*, 7:151103–151119, 2019.

[29] Diego Gragnaniello, Francesco Marra, Giovanni Poggi, and Luisa Verdoliva. Perceptual quality-preserving black-box attack against deep learning image classifiers. *arXiv preprint arXiv:1902.07776*, 2019.

[30] Andras Rozsa, Ethan M Rudd, and Terrance E Boult. Adversarial diversity and hard positive generation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2016.

[31] Fei Gao, Dacheng Tao, Xinbo Gao, and Xuelong Li. Learning to rank for blind image quality assessment. *IEEE Transactions on Neural Networks and Learning Systems*, 26(10):2275–2290, 2015.

[32] Lin Ma, Long Xu, Yichi Zhang, Yihua Yan, and King Ngi Ngan. No-reference retargeted image quality assessment based on pairwise rank learning. *IEEE Transactions on Multimedia*, 18(11):2228–2237, 2016.

[33] Zhengli Zhao, Dheeru Dua, and Sameer Singh. Generating natural adversarial examples. In *Proc. International Conference on Learning Representations*, 2017.

[34] Hongjun Wang, Guangrun Wang, Ya Li, Dongyu Zhang, and Liang Lin. Transferable, controllable, and inconspicuous adversarial attacks on person re-identification with deep mis-ranking. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2020.

[35] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *Proc. International Conference on Learning Representations*, 2014.

[36] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.

[37] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[38] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press Cambridge, 1998.

[39] Wenjie Ruan, Xiaowei Huang, and Marta Kwiatkowska. Reachability analysis of deep neural networks with provable guarantees. In *Proc. International Joint Conference on Artificial Intelligence*, 2018.

[40] Bolin Gao and Lacra Pavel. On the properties of the softmax function with application in game theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.

[41] Francesco Orabona. Almost sure convergence of sgd on smooth non-convex functions. https://parameterfree.com/2020/10/05/almost-sure-convergence-of-sgd-on-smooth-non-convex-functions, 2020.

[42] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.

[43] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[44] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016.

[45] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proc. International Conference on Learning Representations*, 2015.

[46] Zhengyu Zhao, Zhuoran Liu, and Martha Larson. Towards large yet imperceptible adversarial image perturbations with perceptual color distance. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pages 1039–1048, 2020.

[47] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[48] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In *Proc. International Conference on Learning Representations*, 2018.