# Low Frequency Adversarial Perturbation

**Chuan Guo, Jared S. Frank, Kilian Q. Weinberger**
Cornell University, Ithaca, NY 14853, USA
{cg563, jsf239, kqw4}@cornell.edu

## Abstract

Adversarial images aim to change a target model's decision by minimally perturbing a target image. In the black-box setting, the absence of gradient information often renders this search problem costly in terms of query complexity. In this paper we propose to restrict the search for adversarial images to a low frequency domain. This approach is readily compatible with many existing black-box attack frameworks and consistently reduces their query cost by 2 to 4 times. Further, we can circumvent image transformation defenses even when both the model and the defense strategy are unknown. Finally, we demonstrate the efficacy of this technique by fooling the Google Cloud Vision platform with an unprecedented low number of model queries.

## 1 INTRODUCTION

As machine learning models enjoy widespread adoption, their security becomes a relevant topic for consideration. Recent studies have shown that existing methods lack robustness against imperceptible changes to the input [Biggio et al., 2013; Szegedy et al., 2014], and many deployed computer vision and speech recognition systems have been compromised [Carlini and Wagner, 2018; Cisse et al., 2017; Ilyas et al., 2018; Liu et al., 2016; Melis et al., 2017]. This presents a realistic security threat in critical applications such as autonomous driving, where an adversary may manipulate road signs to cause control malfunction while remaining hidden to the naked eye [Evtimov et al., 2017].

Most existing attack algorithms, both white-box [Carlini and Wagner, 2017; Moosavi-Dezfooli, Fawzi, and Frossard, 2016; Szegedy et al., 2014] and black-box



Figure 1: A sample low frequency adversarial image produced by black-box attack.

[Brendel, Rauber, and Bethge, 2017; Chen et al., 2017; Ilyas et al., 2018; Tu et al., 2018], function by searching the full space of possible perturbations to find noise patterns that alter the behavior of convolutional filters. In this high dimensional space many solutions exist and search algorithms tend to almost exclusively result in high frequency solutions, *i.e.* small pixel-wise perturbations dispersed across an image. White-box attacks can be guided by gradient information and tend to have low query complexity (as low as 10 gradients on ResNet/ImageNet).In contrast, black-box attacks do not enjoy such benefits. For example, the search for successful ResNet/ImageNet attacks still requires on the order of $10^4 - 10^5$ queries.

Motivated by these shortcomings, we propose a radical departure from the existing, high-frequency adversarial perturbation attacks and we explicitly restrict the search space of adversarial directions to the low frequency subspace. Constructing low frequency adversarial perturbation has several advantages: As black-box attacks generally require random sampling in the image space, its high-dimensionality causes the attack algorithm to sample many non-adversarial directions, resulting in a high query complexity on the order of the image dimensionality. In the low frequency subspace adversarial directions may occur in much higher density – lowering query complexity significantly. Moreover, many successful defenses against black-box attacks rely on removing high frequency signal with a low-pass filter [Dziugaite, Ghahramani, and Roy, 2016; Guo et al., 2017; Xu, Evans, and Qi, 2017], and operating in low frequency space promises to bypass these

image transformation defenses.

In this paper we show that adversarial perturbations do indeed exist abundantly in a very low-dimensional low frequency subspace. We demonstrate that two popular black-box attacks – the boundary attack [Brendel, Rauber, and Bethge, 2017] and the natural evolution strategies (NES) attack [Ilyas et al., 2018] – can be readily restricted to such a low frequency domain. Figure 1 shows a sample black-box adversarial image with low frequency perturbation produced by the boundary attack. Our experiments demonstrate that a dimensionality reduction to a mere $1/64$ of the original space still yields near-optimal adversarial perturbations. Experimental results confirm our conjectured benefits in the black-box setting:

1. The boundary attack with low frequency perturbation requires dramatically fewer model queries to find an adversarial image. The modified attack produces adversarial images with imperceptible change on ImageNet (ResNet-50 [He et al., 2016]) after approximately 1000 median number of model queries – a 4x reduction compared to vanilla boundary attack.

2. The NES attack enjoys significant improvement of query efficiency with this simple modification, resulting in a consistent 2x speed-up across all images. The median number of queries required for a *targeted* black-box attack using low frequency NES is only around $12,000$.

3. Using low frequency perturbation circumvents image transformation defenses such as JPEG compression [Dziugaite, Ghahramani, and Roy, 2016] and bit depth reduction [Xu, Evans, and Qi, 2017], which have not exhibited vulnerability to black-box attacks prior to our work.

4. Finally, we employ the low frequency boundary attack to fool the Google Cloud Vision platform with an unprecedented 1000 model queries — demonstrating its cost effectiveness and real world applicability.

## 2 BACKGROUND

In the study of adversarial examples in image classification, the goal of an attacker is to alter the model's prediction by adding an imperceptible perturbation to a natural image. Formally, for a given classification model $h$ and an image $\mathbf{x}$ on which the model correctly predicts $y = h(\mathbf{x})$, the adversary aims to find a perturbed image $\mathbf{x}'$ that solves the following constrained optimization problem:

$$\min_{\mathbf{x}'} \ \delta(\mathbf{x}, \mathbf{x}') \text{ subject to } h(\mathbf{x}') \neq y.$$

The function $\delta$ measures the perceptual difference between the original and adversarial images, and is often approximated by mean squared error (MSE), the Euclidean

norm $\| \cdot \|_2$ or the max-norm $\| \cdot \|_\infty$. An attack is considered successful if the perturbed image is imperceptibly different, i.e., $\delta(\mathbf{x}, \mathbf{x}') \leq \rho$ for some small $\rho > 0$. This attack goal defines an *untargeted attack*, since the attack goal is to alter the prediction on the perturbed image to any incorrect class $h(\mathbf{x}') \neq y$. In contrast, a *targeted* attack aims to produce perturbed images that the model predicts as some pre-specified target class.

When constructing adversarial images, the attacker may have various degrees of knowledge about the model $h$, including the training data and/or procedure, model architecture, or even all of its parameters. The attack may also adaptively query $h$ on chosen inputs before producing the adversarial images and obtain gradients from $h$. These different threat models can be roughly categorized into *white-box*, where the attacker has full knowledge about $h$ and how it is trained, or *black-box*, where the attacker can only query $h$, and has limited knowledge about its architecture or training procedure.

**White-box attacks.** When given access to the model entirely, the adversary may minimize the correct class prediction probability directly to cause misclassification [Carlini and Wagner, 2017; Goodfellow, Shlens, and Szegedy, 2015; Kurakin, Goodfellow, and Bengio, 2016; Madry et al., 2017]. For a given input $\mathbf{x}$ and correct class $y$, the adversary defines a loss function $\ell_y(\mathbf{x}')$ so that the loss value is low when $h(\mathbf{x}') \neq y$. One example of such $\ell$ is the margin loss

$$\ell_y(\mathbf{x}') = \max\left( Z(\mathbf{x}')_y - \max_{y' \neq y} Z(\mathbf{x}')_{y'} + \kappa, 0 \right) \quad (1)$$

used in [Carlini and Wagner, 2017], where $Z$ is the logit output of the network. The loss diminishes to zero only if the logit of at least one class exceeds that of the correct class, $y$, by $\kappa$ or more. The adversary can then solve

$$\min_{\mathbf{x}'} \ell_y(\mathbf{x}') + \lambda \delta(\mathbf{x}, \mathbf{x}')$$

with a suitable hyperparameter $\lambda$ to constrain the perturbation to be small while ensuring misclassification.

**Black-box attacks.** In certain scenarios, the white-box threat model does not reflect the true capability of an attacker. For example, when attacking machine learning services such as Google Cloud Vision, the attacker only has access to a limited number of function calls against the target model, and does not have knowledge about the training data. *Transfer-based attacks* [Liu et al., 2016; Papernot et al., 2017; Tramèr et al., 2017] utilize a substitute model that the attacker trains to imitate the target model, and constructs adversarial examples on the substitute model using white-box attacks. For this attack to succeed, the target model must be similar to the substitute model and is trained on the same data distribution.

*Gradient estimation attacks* use techniques such as finite difference [Chen et al., 2017; Tu et al., 2018] and natural evolution strategies [Ilyas et al., 2018; Ilyas, Engstrom, and Madry, 2018] to estimate the gradient from input-output pairs, thus enabling gradient-based white-box attacks. This type of attack requires the model to output class scores or probabilities, and generally requires a number of model queries proportional to the image size. In contrast, *decision-based attacks* [Brendel, Rauber, and Bethge, 2017; Ilyas et al., 2018] utilize only the discrete classification decisions from a model and is applicable in all scenarios, but is generally more difficult to execute.

# 3 LOW FREQUENCY IMAGE SUBSPACE

The inherent query inefficiency of gradient estimation and decision-based attacks stems from the need to search over or randomly sample from the high-dimensional image space. Thus, their query complexity depends on the relative adversarial subspace dimensionality compared to the full image space. One way to improve these methods is to find a low-dimensional subspace that contains a high density of adversarial examples, which enables more efficient sampling of useful attack directions.

Methods in image compression, in particular the celebrated JPEG codec [Wallace, 1991], have long exploited the fact that most of the critical content-defining information in natural images live in the low end of the frequency spectrum, whereas high frequency signal is often associated with noise. It is therefore plausible to assume that CNNs are trained to respond especially to low-frequency patterns in order to extract class-specific signatures from images. Hence, we propose to target CNN based approaches by restricting the search space for adversarial directions to the low-frequency spectrum – essentially targeting these class defining signatures directly.

**Discrete cosine transform.** The JPEG codec utilizes the *discrete cosine transform* (DCT), which decomposes a signal into cosine wave components, to represent a natural image in frequency space. More precisely, given a 2D image $X \in \mathbb{R}^{d \times d}$, define basis functions

$$\phi_d(i,j) = \cos\left[\frac{\pi}{d}\left(i + \frac{1}{2}\right)j\right]$$

for $1 \leq i, j \leq d$. The DCT transform $V = \mathrm{DCT}(X)$ is:

$$V_{j_1,j_2} = N_{j_1} N_{j_2} \sum_{i_1=0}^{d-1} \sum_{i_2=0}^{d-1} X_{i_1,i_2} \phi_d(i_1,j_1)\phi_d(i_2,j_2),$$

where $N_j = \sqrt{\frac{1}{d}}$ if $j = 0$ and $N_j = \sqrt{\frac{2}{d}}$ otherwise. Here, $N_{j_1}, N_{j_2}$ are normalization terms included to ensure the
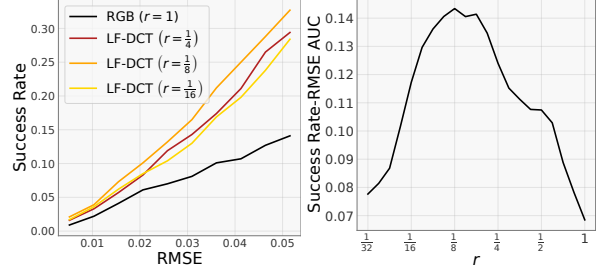


Figure 2: (Left) Comparison of (attack) success rate after perturbation by random spherical noise in RGB vs. LF-DCT space. Using low frequency noise improves the success rate dramatically. (Right) Area under the *success rate-$\rho$* curve. The highest density of adversarial images appears to lie around $r = 1/8$.

transformation is isometric, i.e. $\|X\|_2 = \|\mathrm{DCT}(X)\|_2$. The entry $V_{i,j}$ corresponds to the magnitude of wave $\phi_d(i,j)$, with lower frequencies represented by lower $i, j$. Further, DCT is invertible, with inverse $X = \mathrm{IDCT}(V)$,

$$X_{i_1,i_2} = \sum_{j_1=0}^{d-1} \sum_{j_2=0}^{d-1} N_{j_1} N_{j_2} V_{j_1,j_2} \phi_d(i_1,j_1)\phi_d(i_2,j_2). \tag{2}$$

For images containing multiple color channels, both DCT and IDCT can be applied channel-wise independently.

**Sampling low frequency noise.** In order to facilitate efficient search for attack directions in low frequency space, we need to be able to sample random perturbations confined to this subspace. We can achieve this with the inverse DCT transform by considering the top-left $rd \times rd$ entries of $V$ for some ratio parameter $r \in (0, 1]$. These coefficients correspond to cosine waves with long periods, hence low frequency. Given any distribution $\mathcal{D}$ (e.g. uniform, or Gaussian) over $\mathbb{R}^{d \times d}$, we can sample a random matrix $\tilde{\eta} \in \mathbb{R}^{d \times d}$ in frequency space so that

$$\tilde{\eta}_{i,j} = \begin{cases} x_{i,j} \sim \mathcal{D} & \text{if } 1 \leq i, j \leq rd \\ 0 & \text{otherwise.} \end{cases}$$

Using the inverse DCT mapping, the corresponding noise "image" in pixel space is defined by $\eta = \mathrm{IDCT}(\tilde{\eta})$. By definition, $\eta$ has non-zero cosine wave coefficients only in frequencies lower than $rd$. When the pixel space contains multiple color channels, we can sample each channel independently using the same strategy. We denote this distribution of low frequency noise as $\mathrm{IDCT}_r(\mathcal{D})$ and the sub-space as low frequency DCT (LF-DCT) space.

**Low frequency noise success rate.** We postulate that CNNs are more sensitive to changes in the LF-DCT subspace, hence admitting a higher density of adversarial perturbations. To empirically substantiate this hypothesis, we compare the success rate of random noise in RGB vs.

LF-DCT space for a ResNet-50 architecture [He et al., 2016]. We sample the noise vector $\eta$ uniformly from the surface of the unit sphere of radius $\rho > 0$ in the $rd \times rd$ LF-DCT space and project it back to RGB through the IDCT transform. For $r = 1$ this procedure is identical to sampling directly on the surface of a unit sphere in pixel space, as IDCT is a linear, orthonormal transformation.

Figure 2 (left) shows these success rates as a function of the root mean squared error (RMSE = $\frac{\rho}{\sqrt{3d}}$) between the perturbed and original image, averaged over 1000 randomly chosen images from ImageNet [Deng et al., 2009]. Several trends emerge: 1. As expected, the success rate increases with the magnitude of perturbation $\rho$ across all values of $r$; 2. There appears to be a sweet spot around $r = 1/8$, which corresponds to a reduction of dimensionality by $1/64$; 3. The worst success rate is achieved with $r = 1$, which corresponds to no dimensionality reduction (and is identical to sampling in the original RGB space).

To further investigate the dimensionality "sweet spot", the right plot shows the area under the success rate vs. RMSE curve for various values of $r$, within a fixed range of $\rho \in [0, 20]$. Here, a higher value corresponds to a faster increase in success rate with larger perturbation radius. In agreement with the left plot, the optimal frequency ratio is around $r = 1/8$.

### 3.1 Universality of low frequency subspace

Results in the previous section support our hypothesis that restricting the search space to LF-DCT substantially increases the sample success rate of random adversarial directions. However, the dimensionality reduction does impose a restriction on the possible solutions of attack algorithms. To examine the effects of this limitation, we apply our low-frequency restriction to white-box attacks by projecting the gradient onto the LF-DCT space.

**Low frequency gradient descent.** Let $\ell_y$ denote the adversarial loss, e.g. Equation 1. For a given $r \in (0, 1]$ and $v \in \mathbb{R}^{rd \times rd}$, define $V \in \mathbb{R}^{d \times d}$ by

$$V_{i,j} = \begin{cases} v_{i,j} & \text{if } 1 \leq i, j \leq rd \\ 0 & \text{otherwise,} \end{cases}$$

The wave coefficient matrix $V$ contains $v$ as its submatrix and only includes frequencies lower than $rd$. The low frequency perturbation domain can then be parametrized as $\Delta = \text{IDCT}(V)$. To optimize with gradient descent, let $\bar{\Delta}$ and $\bar{V}$ be vectorizations of $\Delta$ and $V$, i.e., $\bar{\Delta}_{i_1 * d + i_2} = \Delta_{i_1, i_2}$ and similarly for $\bar{V}$. From Equation 2, it is easy to see that each coordinate of $\bar{\Delta}$ is a linear function of $\bar{V}$, hence IDCT is a linear transformation, whose adjoint is precisely the linear transformation defined by DCT. For any vector $\mathbf{z}$, its right-product with the Jacobian of IDCT
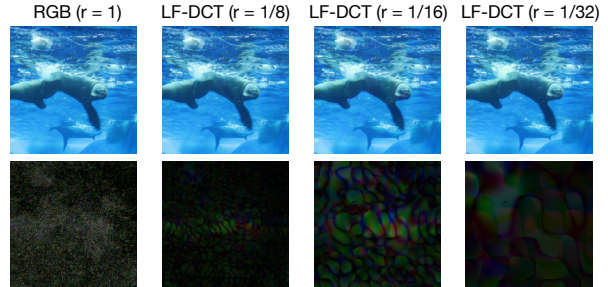


Figure 3: A sample image perturbed by the Carlini-Wagner attack using the full image space and low frequency space with different $r$. The adversarial perturbation (second row) has clearly different pattern across different frequency ranges.

| | $d'$ | MSE | Success Rate (%) |
|---|---|---|---|
| RGB ($r = 1$) | 150528 | $2.78 \times 10^{-5}$ | 100.0 |
| LF-DCT ($r = 1/8$) | 2352 | $6.94 \times 10^{-5}$ | 100.0 |
| LF-DCT ($r = 1/16$) | 588 | $1.61 \times 10^{-4}$ | 95.5 |
| LF-DCT ($r = 1/32$) | 147 | $1.56 \times 10^{-4}$ | 56.0 |

Table 1: Average MSE and accuracy after Carlini-Wagner attack with different frequency ratios $r$. $d' = 3 \times rd \times rd$ is the effective adversarial space dimensionality. At $r = 1/8$, optimizing in the frequency space of dimensionality 2352 is as effective as optimizing in the full image space.

is given by $J_{\text{IDCT}} \cdot \mathbf{z} = \text{DCT}(\mathbf{z})$. Thus we may apply the chain rule to compute

$$\frac{\partial \ell}{\partial V} = \text{DCT}\left(\frac{\partial \ell}{\partial \Delta}\right), \frac{\partial \ell}{\partial v} = \left[\frac{\partial \ell}{\partial V}\right]_{1:rd, 1:rd},$$

which is equivalent to applying DCT to the gradient and dropping the high frequency coefficients.

**Adversarial optimality in low frequency subspace.** Table 1 shows average perturbation MSE and model accuracy after the Carlini-Wagner attack [Carlini and Wagner, 2017] in low frequency space. The original attack in pixel space corresponds to $r = 1$. The images have three color channels and the effective subspace dimensionality is $d' = 3 \times rd \times rd$. For $r = 1/8$, the attack can achieve perfect (100%) success rate, while the resulting MSE is only roughly 3 times larger — despite that the search space dimensionality is only $1/64$ of the full image space. This result further supports that the density of adversarial examples is much higher in the low frequency domain, and that searching exclusively in this restricted subspace consistently yields near-optimal adversarial perturbations. As expected, choosing a very small frequency ratio eventually impacts success rate, as the subspace dimensionality is too low to admit adversarial perturbations. Figure 3 shows the resulting adversarial images and perturbations corresponding to frequency ratios $r$. All perturbations
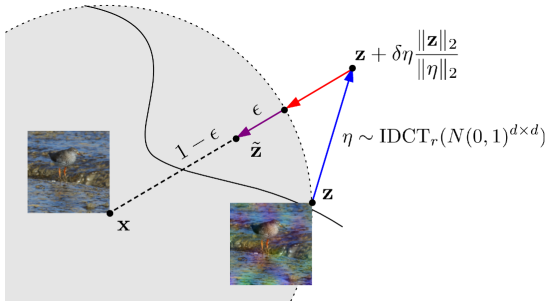
Figure 4: Illustration of a single iteration of the low frequency boundary attack. Instead of sampling the noise matrix $\eta$ from $N(0,1)^{d\times d}$, we sample a low frequency noise matrix by applying IDCT to the Gaussian noise while removing high frequency components.

are imperceptible but when isolated (bottom row) reveal increasingly smooth patterns as r decreases.

**Advantages of low frequency perturbation.** While the remainder of this paper focuses on the benefits of low frequency adversarial perturbation in the black-box setting, we highlight that there are advantages in the white-box setting as well. Sharma, Ding, and Brubaker [2019] showed that low frequency gradient-based attacks enjoy greater efficiency and can transfer significantly better to defended models. In particular, their attack is able to completely circumvent all of the top-placing defense entries at the NeurIPS 2017 competition. Furthermore, they observe that the benefit of low frequency perturbation is not merely due to dimensionality reduction — perturbing exclusively the high frequency components does not give the same benefit.

## 4 APPLICATION TO BLACK-BOX ATTACKS

Many existing black-box attacks proceed by iteratively adding random noise to the current image and evaluating the model to determine the direction to move towards. Given our insights regarding the effectiveness of low frequency perturbations, we propose its use as a universal tool for improving the query efficiency of black-box attacks. We conduct case studies on the boundary attack [Brendel, Rauber, and Bethge, 2017] and the NES attack [Ilyas et al., 2018] to demonstrate the efficacy and accessibility of our method.

### 4.1 Case study: Boundary attack

The boundary attack uses an iterative update rule to gradually move the adversarial image closer to the original image, maintaining that the image remains adversarial

at each step. Starting from random noise, the algorithm samples a noise matrix $\eta \sim N(0,1)^{d\times d}$ at each iteration and adds it to the current iterate $\mathbf{z}$ after appropriate scaling. This point is then projected onto the sphere of center $\mathbf{x}$, the target image, and radius $\|\mathbf{z}\|_2$ so that the next iterate never moves away from $\mathbf{x}$. Finally, we contract towards $\mathbf{x}$ by $\epsilon$, and the new iterate $\tilde{\mathbf{z}}$ is accepted only if it remains adversarial. This guarantees that terminating the algorithm at any point still results in an adversarial image, but the perturbation magnitude reduces with each contraction step.

**Modification.** To construct low frequency perturbation using the boundary attack, we constrain the noise matrix $\eta$ to be sampled from $\mathrm{IDCT}_r(N(0,1)^{d\times d})$ instead. Figure 4 illustrates the modified attack. Sampling low frequency noise instead of Gaussian noise is particularly beneficial to the boundary attack: After adding the noise matrix $\eta$, if the iterate is not adversarial, the algorithm must re-sample a noise matrix and perform another model query. By restricting to the low frequency subspace, which has a larger density of adversarial directions, this step succeeds more often, speeding up convergence towards the target image. We term this variant of the boundary attack as *low frequency boundary attack* (LF-BA) and the original boundary attack as RGB-BA.

**Hyperparameters.** The boundary attack has two hyperparameters: noise step size $\delta$ and contraction step size $\epsilon$. Both step sizes are adjusted based on the success rate of the past few candidates, *i.e.*, if $\tilde{\mathbf{z}}$ is accepted often, we can contract towards the target $\mathbf{x}$ more aggressively by increasing $\epsilon$ and vice versa, and similarly for $\delta$. For the low frequency variant, we find that fixing $\delta$ to a large value is beneficial for speeding up convergence, while also reducing the number of model queries by half. For all experiments, we fix $\delta = 0.2$ and initialize $\epsilon = 0.01$.

Selecting the frequency ratio $r$ is more crucial. Different images may admit adversarial perturbations at different frequency ranges, and thus we would like the algorithm to automatically discover the right frequency on a per-image basis. We use Hyperband [Li et al., 2016], a bandit-type algorithm for selecting hyperparameters, to optimize the frequency ratio $r$. We initialize Hyperband with multiple runs of the attack for every frequency ratio $r \in \{\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}\}$. Repeatedly after $T'$ iterations, the least successful half of the parallel runs is terminated until one final frequency remains. This setting is continued until the total number of model queries reaches $T$.

### 4.2 Case study: NES attack

Natural evolution strategies (NES) [Wierstra et al., 2014] is a black-box optimization technique that has been re-
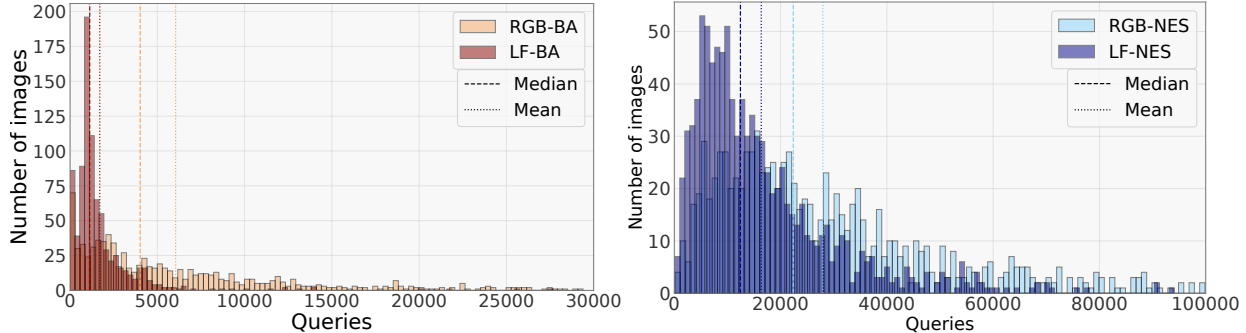
Figure 5: Distribution of the number of queries required for a successful attack (defined as achieving a perturbation MSE of 0.001 or lower for RGB-BA/LF-BA). See text for details.

cently proposed for its use in black-box attacks [Ilyas et al., 2018]. The attacker constructs the adversarial image $\mathbf{z}$ by minimizing a continuous-valued adversarial loss $\ell$ returned by black-box query to the model. However, instead of minimizing $\ell$ directly, the NES attack minimizes the loss at all points near $\mathbf{z}$. More precisely, we specify a search distribution $\mathcal{D}$ and minimize:

$$\min_{\mathbf{z}} \mathbb{E}_{\eta \sim \mathcal{D}}[\ell(\mathbf{z} + \eta)] \text{ subject to } d(\mathbf{x}, \mathbf{z}) \leq \rho, \quad (3)$$

where $\rho$ is some perceptibility threshold. When the search distribution $\mathcal{D}$ is chosen to be an isotropic Gaussian, i.e. $\mathcal{D} = N(0, \sigma^2)^{d \times d}$, the gradient of the objective function in Equation 3 becomes

$$\nabla_{\mathbf{z}} \mathbb{E}_{\eta \sim \mathcal{D}}[\ell(\mathbf{z} + \eta)] = \frac{1}{\sigma^2} \mathbb{E}_{\eta \sim \mathcal{D}}[\ell(\mathbf{z} + \eta) \cdot \eta].$$

Thus, Equation 3 can be minimized with stochastic gradient descent by sampling a batch of noise vectors $\eta_1, \ldots, \eta_m \sim N(0, \sigma^2)^{d \times d}$ and computing the (mini-batch) stochastic gradient

$$\nabla_{\mathbf{z}} \mathbb{E}_{\eta \sim \mathcal{D}}[\ell(\mathbf{z} + \eta)] \approx \frac{1}{m\sigma^2} \sum_{i=1}^{m} \ell(\mathbf{z} + \eta_i) \cdot \eta_i. \quad (4)$$

One way to interpret this update rule is that the procedure pushes $\mathbf{z}$ away from regions of low adversarial density — directions $\eta_i$ for which $\ell(\mathbf{z}+\eta_i)$ is high. The perceptibility constraint can be enforced by projecting to the feasible region at every step. For this attack, the max-norm $\|\cdot\|_{\infty}$ is used as the perceptibility metric, hence the projection step reduces to clipping of each dimension in the adversarial perturbation to the range $[-\rho, \rho]$.

**Modification.** The low frequency distribution defined in section 3 can be readily incorporated into the NES attack. We replace the Gaussian search distribution with its low frequency version, i.e. we sample a batch of noise vectors $\eta_1, \ldots, \eta_m \sim \text{IDCT}_r(N(0, \sigma^2)^{d \times d})$ instead. The stochastic gradient remains identical to Equation 4. Note that since each $\eta_i$ is low-frequency, this process results

in a low frequency adversarial perturbation. We term the original NES attack using search distribution in pixel space as RGB-NES and the low frequency variant as *low frequency NES* (LF-NES).

**Hyperparameters.** The NES attack has two hyperparameters: $\rho$, which controls the perceptibility of adversarial perturbation, and $\sigma$, which controls the width of the search distribution. We set $\rho = 0.03$ to match the average $L_2$-norm of perturbations generated by RGB-BA/LF-BA, and set $\sigma = 0.001$ as suggested by the authors. Intriguingly, the frequency ratio $r$ is not very sensitive for LF-NES. Setting a single value of $r$ for all images is sufficiently effective, and we choose the same value of $r = 1/2$ in all of our experiments for simplicity.

## 5 EMPIRICAL EVALUATION

We empirically validate our claims that black-box attacks in low frequency space possess the aforementioned desirable properties. For all experiments, we use the default PyTorch pretrained ResNet-50 model for RGB-BA/LF-BA and the Tensorflow-Slim pretrained ResNet-50 model[1] for RGB-NES/LF-NES.

Both RGB-BA and LF-BA use a 10 step binary search along the line joining the random initialization and the target image before starting the attack. Our implementation of the boundary attack in PyTorch has comparable performance to the official implementation by Brendel, Rauber, and Bethge [2017] while being significantly faster. We use the official implementation of NES in Tensorflow and modify it to use low frequency search distribution. We release our code[2] publicly for reproducibility.

**Settings.** For experiments on ImageNet [Deng et al., 2009], we evaluate both untargeted attack (RGB-BA/LF-

[1]https://github.com/tensorflow/models/tree/master/research/slim
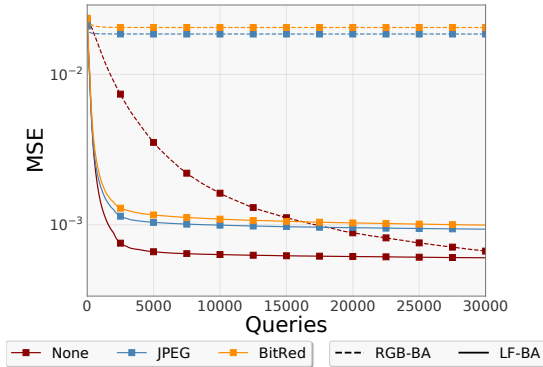[2]https://github.com/cg563/low-frequency-adversarial

Figure 6: Average (log) MSE across queries for RGB-BA and LF-BA against different image transformation defenses. Against JPEG and bit depth reduction defenses, RGB-BA fail to make progress. In contrast, LF-BA can successfully circumvent these defenses and reduce average MSE to 0.001 after $30,000$ model queries.



Figure 7: Plot of the average perturbation MSE across iterations for LF-BA using different frequency ratios $r$. See text for details.

BA) and targeted attack (RGB-NES/LF-NES) to a random class against a pretrained ResNet-50 [He et al., 2016] model. Each test image is randomly selected from the validation set while ensuring correct prediction by the respective models. For RGB-BA/LF-BA, the adversary can only access the binary output of the model corresponding to whether the input is classified as the original label. For RGB-NES/LF-NES, the adversary can obtain the cross entropy loss of the model against the target label.

We limit the attack algorithm to $30,000$ queries for untargeted attack, corresponding to $30,000$ iterations for LF-BA and $15,000$ for RGB-BA[3], and $100,000$ queries for targeted attack using RGB-NES/LF-NES. For LF-BA, we select the frequency ratio $r$ using Hyperband by halving the number of parallel runs every $T' = 500$ iterations.

**Query histogram.** Figure 5 shows the histogram of the number of model queries required for a successful attack over 1000 sampled images. The left plot shows result for untargeted attack using RGB-BA/LF-BA. Since the boundary attack maintains an incorrectly labeled image throughout optimization while gradually reducing the perturbation norm, we define success as achieving a sufficiently low perceptibility of $< 0.001$ MSE (or equivalently, an $L_2$-norm of 12.27). The results for targeted attacks using RGB-NES/LF-NES are in the right plot. Only successful runs are included in this plot. We make several key observations:

1. The query distribution of RGB-BA (light orange) and RGB-NES (light blue) are heavy-tailed, that is, the entire range of allowed number of queries is covered, which shows that a large number of model queries is necessary

---
[3]RGB-BA requires two model queries per iteration, one after the noise step and one after the contraction.
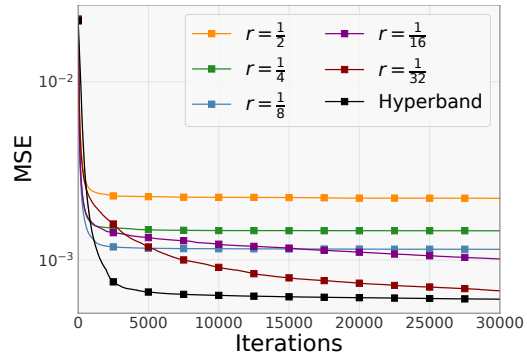
for many images.

2. The histograms of LF-BA (dark red) and LF-NES (dark blue) are shifted left compared to their Gaussian-based counterparts. This demonstrates that using the low frequency noise samples consistently improves the query efficiency of the boundary attack and the NES attack. This effect is especially dramatic for LF-BA, where a large fraction of images require only roughly 1000 model queries to construct.

3. Both the median (dashed line) and mean (dotted line) query counts are significantly reduced when using LF-BA and LF-NES. In particular, LF-BA requires $1128$ median queries, an almost 4x reduction compared to the $4020$ median queries of RGB-BA. Similarly, LF-NES requires $12,444$ median queries, an approximately 2x reduction from the $22,389$ median queries of RGB-NES.

**Selecting frequency ratio $r$.** In Figure 7 we analyze the effect of selecting the hyperparameter $r$ by either fixing it to a pre-defined value or by using Hyperband. The average MSE of adversarial perturbations constructed by LF-BA is plotted against the number of iterations. All averages are computed over the same 1000 random images from the ImageNet validation set.

At higher values of $r$, the perturbation MSE drops rapidly for the first 2500 iterations but progress stalls later on. Lower values of $r$ (e.g. purple and dark red lines) allow the attack algorithm to (relatively) slowly but eventually find an adversarial perturbation with low MSE. This plot demonstrates the need for selecting $r$ adaptively based on the image, as Hyperband (black line) selects the optimal frequency ratio to allow both rapid descent initially and continued progress later on.

**Breaking image transformation defenses.** One common defense strategy against adversarial images is to apply a denoising transformation before feeding it into the model. This style of defense has been shown to be
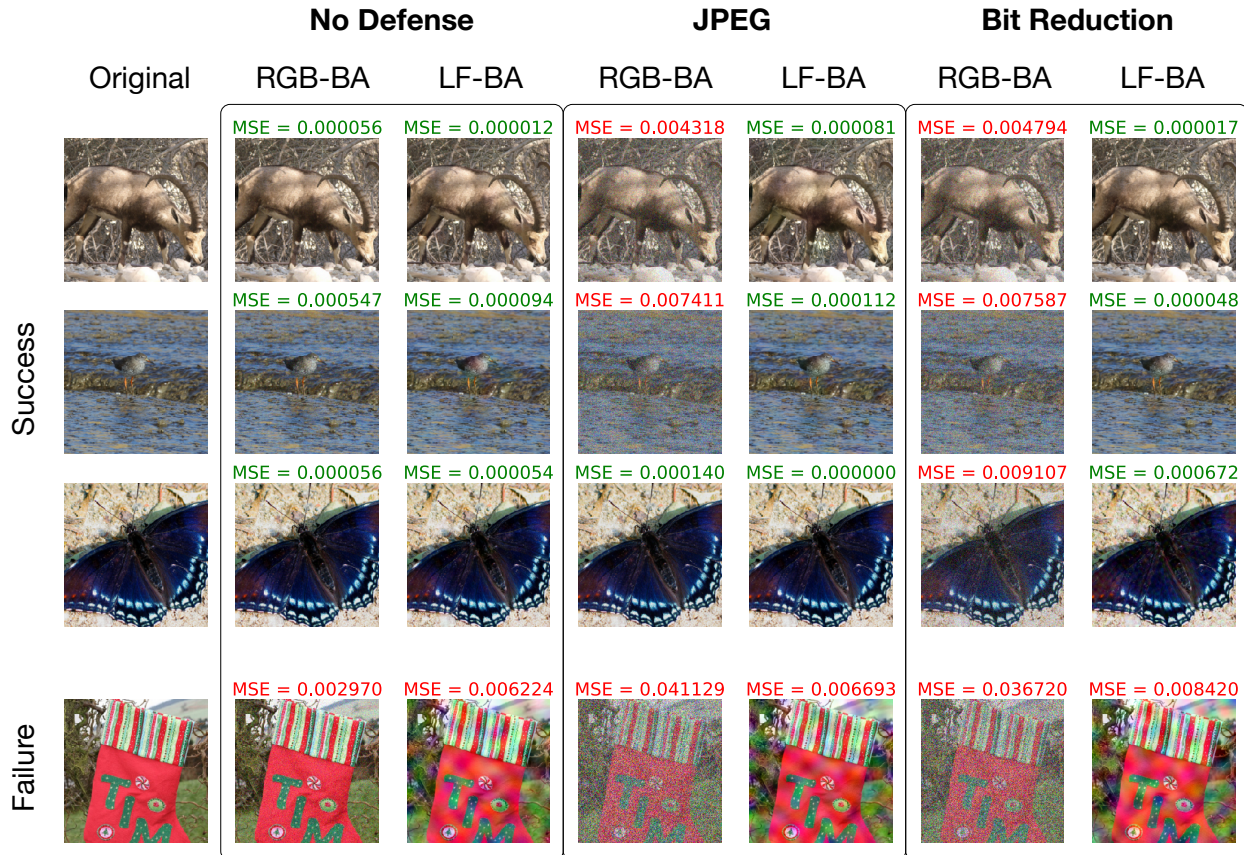
Figure 8: Image samples for attacking image transformation defenses. Perturbation MSE is truncated to 8 decimal places and images with MSE higher than 0.001 are highlighted in red. See text for details.

| | Untargeted | | | | Targeted | | | |
|---|---|---|---|---|---|---|---|---|
| Attack | Average queries | Average $L_2$ | Average MSE | Attack | Average queries | Success rate | Average $L_2$ | Average MSE |
| Opt-attack | 71, 100 | 6.98 | $3.24 \times 10^{-4}$ | AutoZOOM | 13, 525 | 100% | 26.74 | $3.64 \times 10^{-3}$ |
| RGB-BA | 14, 217 | - | - | RGB-NES | 31, 879 | 94.7% | 6.85 | $3.22 \times 10^{-4}$ |
| LF-BA | 2926 | - | - | LF-NES | 17, 558 | 98.6% | 6.92 | $3.18 \times 10^{-4}$ |

Table 2: Comparison of aggregate statistics on ImageNet. All statistics are averaged over 1000 random validation images. See text for details.

highly effective against transfer-based attacks [Guo et al., 2017] and have not exhibited any weakness against black-box attacks to-date. However, we suspect that low frequency perturbations can circumvent this defense since denoising transformations do not typically operate on the lower frequency spectrum.

We test our hypothesis by evaluating RGB-BA and LF-BA against the following image transformation defenses: JPEG compression [Dziugaite, Ghahramani, and Roy, 2016] at quality level 75 and reducing bit depth [Xu, Evans, and Qi, 2017] to 3 bits. To avoid artificially inflating success rate, we choose initial images to be correctly classified after the defensive transformation is applied.

Figure 6 compares the average perturbation (log) MSE across model queries for both attacks on 1000 random images across iterations. Again, we see that LF-BA (solid line) converges significantly faster than RGB-BA (dashed line) when the model is undefended (dark red lines). In fact, it reaches the same average MSE achieved by RGB-BA after 30, 000 model queries in less than 3000 queries – constituting an *order of magnitude* reduction. When either the JPEG (blue lines) or bit depth (orange lines) reduction transformation is applied, RGB-BA fails to make any progress. This result shows that image transformation defenses are very potent against black-box attacks. On the other hand, LF-BA can circumvent these defenses consistently and reduce the average perturbation MSE to approximately 0.001 after 30, 000 model queries. The success of LF-BA *does not* rely on the knowledge of the exact transformation being applied.
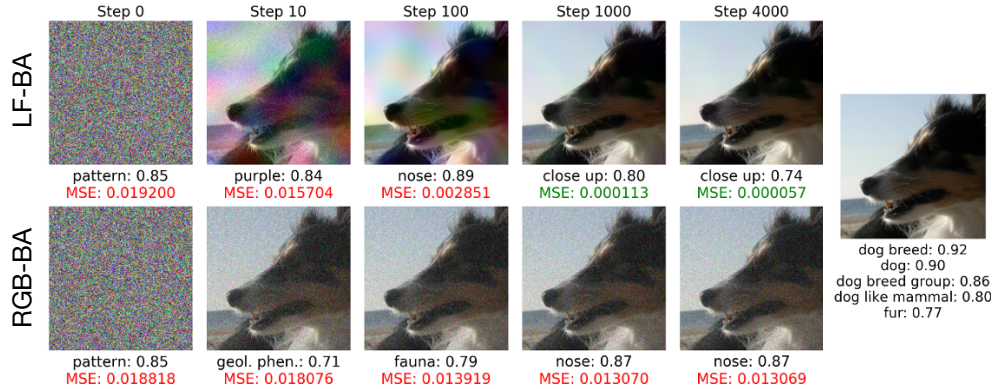
Figure 9: Attacking Google Cloud Vision. MSE of value higher than 0.001 is colored in red. See text for details.

Figure 8 shows adversarially perturbed images produced when attacking different image transformation defenses. On the undefended model, there is no visible difference between the clean image and the perturbed image when attacking with either Gaussian or low frequency noise. On defended models, RGB-BA consistently fails to produce an imperceptible perturbation, while LF-BA is successful with high probability. Note the color patch pattern produced by LF-BA has varying frequency, which is optimally selected by Hyperband. The last image represents a failure case for both RGB-BA and LF-BA.

**Additional baselines.** Table 2 shows aggregate query and perturbation norm statistics for untargeted RGB-BA/LF-BA attacks and targeted RGB-NES/LF-NES attacks in comparison to two additional baselines: Opt-attack [Cheng et al., 2018] and AutoZOOM [Tu et al., 2018]. We duplicate relevant numbers reported in the original paper for both baselines. Since RGB-BA/LF-BA gradually reduce perturbation magnitude at the expense of additional queries, we set a target average $L_2$-norm equal to that of Opt-attack and compare query cost. For RGB-NES/LF-NES, we fix the same maximum number of queries to $100,000$ as AutoZOOM, and compare query count and perturbation magnitude at initial success.

Note that LF-BA requires 5x fewer queries than RGB-BA and 24x fewer queries than Opt-attack to reach the same average $L_2$-norm/MSE, constituting an *order of magnitude* reduction. For targeted attack, LF-NES requires only half as many queries as RGB-NES to reach the same perturbation norm while having higher success rate. Compared to AutoZOOM, LF-NES requires approximately the same number of average queries while achieving a nearly 4x reduction of perturbation $L_2$-norm.

**Attacking Google Cloud Vision.** To demonstrate the realistic threat of low frequency perturbations, we attack Google Cloud Vision, a popular online machine learning service. The platform provides a top concept label-

ing functionality: when given an image, it outputs a list of top (predicted) concepts contained in the image and their associated confidence. We define a successful attack as replacing the formerly highest ranked concept with a new concept that was previously not present in the list, while obtaining an MSE $\leq 0.001$. Figure 9 shows the progression of the boundary attack with Gaussian and low frequency noise across iterations. On the image with original top concept *dog breed*, LF-BA produces an adversarial image with imperceptible difference while changing the top concept to *close-up*. Even with only 1000 model queries, the adversarial perturbation is already reasonably unobtrusive. In contrast, RGB-BA could not find a sufficiently minimal perturbation within 4000 iterations (=8000 queries). Note that neither method makes use of the prediction confidence or the rank of concepts other than the top-1, contrasting with the previous known attack against this platform [Ilyas et al., 2018].

# 6 DISCUSSION AND FUTURE WORK

We have shown that adversarial attacks on images can be performed by exclusively perturbing low frequency portions of the input signal. This approach provides substantial benefits for attacks in the black-box setting and can be readily incorporated into many existing algorithms. Our follow-up work [Guo et al., 2019] that achieves state-of-the-art query efficiency using a simple coordinate descent-style attack also leverages the abundance of adversarial perturbations in the low frequency subspace.

Focusing on low frequency signal is by no means exclusively applicable to images. It is likely that similar approaches can be used to attack speech recognition systems [Carlini and Wagner, 2018] or time series data. Another promising future direction is to find other subspaces that may admit a higher density of adversarial perturbations. Any success in this direction can also provide us with insight into the space of adversarial examples.

# References

Biggio, B.; Corona, I.; Maiorca, D.; Nelson, B.; Šrndić, N.; Laskov, P.; Giacinto, G.; and Roli, F. 2013. Evasion attacks against machine learning at test time. In *Proc. ECML*, 387–402.

Brendel, W.; Rauber, J.; and Bethge, M. 2017. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *CoRR* abs/1712.04248.

Carlini, N., and Wagner, D. A. 2017. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 39–57.

Carlini, N., and Wagner, D. A. 2018. Audio adversarial examples: Targeted attacks on speech-to-text. *CoRR* abs/1801.01944.

Chen, P.; Zhang, H.; Sharma, Y.; Yi, J.; and Hsieh, C. 2017. ZOO: zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec@CCS 2017, Dallas, TX, USA, November 3, 2017*, 15–26.

Cheng, M.; Le, T.; Chen, P.; Yi, J.; Zhang, H.; and Hsieh, C. 2018. Query-efficient hard-label black-box attack: An optimization-based approach. *CoRR* abs/1807.04457.

Cisse, M.; Adi, Y.; Neverova, N.; and Keshet, J. 2017. Houdini: Fooling deep structured prediction models. *CoRR* abs/1707.05373.

Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, 248–255. IEEE.

Dziugaite, G. K.; Ghahramani, Z.; and Roy, D. M. 2016. A study of the effect of JPG compression on adversarial images. *CoRR* abs/1608.00853.

Evtimov, I.; Eykholt, K.; Fernandes, E.; Kohno, T.; Li, B.; Prakash, A.; Rahmati, A.; and Song, D. 2017. Robust physical-world attacks on machine learning models. *CoRR* abs/1707.08945.

Goodfellow, I.; Shlens, J.; and Szegedy, C. 2015. Explaining and harnessing adversarial examples. In *Proc. ICLR*.

Guo, C.; Rana, M.; Cissé, M.; and van der Maaten, L. 2017. Countering adversarial images using input transformations. *CoRR* abs/1711.00117.

Guo, C.; Gardner, J. R.; You, Y.; Wilson, A. G.; and Weinberger, K. Q. 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*, 2484–2493.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proc. CVPR*, 770–778.

Ilyas, A.; Engstrom, L.; Athalye, A.; and Lin, J. 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*, 2142–2151.

Ilyas, A.; Engstrom, L.; and Madry, A. 2018. Prior convictions: Black-box adversarial attacks with bandits and priors. *CoRR* abs/1807.07978.

Kurakin, A.; Goodfellow, I.; and Bengio, S. 2016. Adversarial machine learning at scale. *CoRR* abs/1611.01236.

Li, L.; Jamieson, K. G.; DeSalvo, G.; Rostamizadeh, A.; and Talwalkar, A. 2016. Efficient hyperparameter optimization and infinitely many armed bandits. *CoRR* abs/1603.06560.

Liu, Y.; Chen, X.; Liu, C.; and Song, D. 2016. Delving into transferable adversarial examples and black-box attacks. *CoRR* abs/1611.02770.

Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *CoRR* abs/1706.06083.

Melis, M.; Demontis, A.; Biggio, B.; Brown, G.; Fumera, G.; and Roli, F. 2017. Is deep learning safe for robot vision? adversarial examples against the icub humanoid. *CoRR* abs/1708.06939.

Moosavi-Dezfooli, S.; Fawzi, A.; and Frossard, P. 2016. Deepfool: A simple and accurate method to fool deep neural networks. In *Proc. CVPR*, 2574–2582.

Papernot, N.; McDaniel, P. D.; Goodfellow, I. J.; Jha, S.; Celik, Z. B.; and Swami, A. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab Emirates, April 2-6, 2017*, 506–519.

Sharma, Y.; Ding, G. W.; and Brubaker, M. 2019. On the effectiveness of low frequency perturbations. *CoRR* abs/1903.00073.

Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2014. Intriguing properties of neural networks. In *In Proc. ICLR*.

Tramèr, F.; Kurakin, A.; Papernot, N.; Boneh, D.; and McDaniel, P. D. 2017. Ensemble adversarial training: Attacks and defenses. *CoRR* abs/1705.07204.

Tu, C.; Ting, P.; Chen, P.; Liu, S.; Zhang, H.; Yi, J.; Hsieh, C.; and Cheng, S. 2018. Autozoom: Autoencoder-based zeroth order optimization method for attacking black-box neural networks. *CoRR* abs/1805.11770.

Wallace, G. K. 1991. The jpeg still picture compression standard. *Commun. ACM* 34(4):30–44.

Wierstra, D.; Schaul, T.; Glasmachers, T.; Sun, Y.; Peters, J.; and Schmidhuber, J. 2014. Natural evolution strategies. *Journal of Machine Learning Research* 15(1):949–980.

Xu, W.; Evans, D.; and Qi, Y. 2017. Feature squeezing: Detecting adversarial examples in deep neural networks. *CoRR* abs/1704.01155.