# GreedyFool: Multi-Factor Imperceptibility and Its Application to Designing a Black-box Adversarial Attack

Hui Liu   Bo Zhao   Minzhi Ji
School of Cyber Science and Engineering, Wuhan University
Wuhan, 430072 China
{liuh824, zhaobo, jiminzhi}@whu.edu.cn

Peng Liu
College of Information Sciences and Technology Pennsylvania State University
PA, 16801 US
pliu@ist.psu.edu

## Abstract

*Adversarial examples are well-designed input samples, in which perturbations are imperceptible to the human eyes, but easily mislead the output of deep neural networks (DNNs). Existing works synthesize adversarial examples by leveraging simple metrics to penalize perturbations, that lack sufficient consideration of the human visual system (HVS), which produces noticeable artifacts. To explore why the perturbations are visible, this paper summarizes four primary factors affecting the perceptibility of human eyes. Based on this investigation, we design a multi-factor metric* MulFactorLoss *for measuring the perceptual loss between benign examples and adversarial ones. In order to test the imperceptibility of the multi-factor metric, we propose a novel black-box adversarial attack that is referred to as GreedyFool. GreedyFool applies differential evolution to evaluate the effects of perturbed pixels on the confidence of a target DNN, and introduces greedy approximation to automatically generate adversarial perturbations. We conduct extensive experiments on the ImageNet and CIFRA-10 datasets and a comprehensive user study with 60 participants. The experimental results demonstrate that* MulFactorLoss *is a more imperceptible metric than the existing pixelwise metrics, and GreedyFool achieves a 100% success rate in a black-box manner.*

## 1. Introduction

Existence of adversarial examples [2, 4, 17, 18, 20, 23, 25, 32, 37, 40] is one of the most important findings in deep learning research [7, 13, 26, 36, 39]. Since adversarial examples could enable researchers to gain profound understanding about the nature of deep learning models and how to safely deploy them in real-world AI systems, they have attracted a great amount of interests in the research community in recent years. Based on a widely-accepted notion of adversarial examples, in which certain changes to a sample of input data (e.g. an image) are (largely) **imperceptible** to the human eye, but cause a deep learning model to **misclassify** (or mis-predict) the input, researchers have been answering two basic research questions: (1) Why could certain slight changes trigger a deep learning model to misclassify? (2) Why could the changes needed to "fool" a deep learning model be imperceptible to the human eye? Besides answering the two basic research questions, researchers have also been developing integrated approaches to generating adversarial examples. By "integrated", we mean that the research findings in answering both questions are leveraged. For example, when a deep neural network (DNN) [12, 14, 15, 30] is used to classify images, generating adversarial examples could be formalized as an optimization problem with imperceptibility (i.e. small perceptual loss) as a main constraint [22, 24].

Although the first research question has been extensively studied in the literature [8, 10, 16, 29, 34], the second question is less full investigation. In fact, recent studies attempted to study the second question from different angles [2, 17, 21, 25, 35, 41]. However, what are the primary factors affecting the perceptibility of human eyes? how do these factors combine to better measure the imperceptibility of perturbations? They does not give a comprehensive answer. In order to clearly show the limitations of the existing works, let's firstly summarize the primary factors affecting the HVS. To our best knowledge, the past (psychophysics) researches on the HVS have identified four primary factors

affecting the **perceptibility** of human eyes: (*F1*) Just noticeable distortion (JND) [35,38]. Human eyes cannot sense a stimulus below the JND, which quantifies the sensitivity of the HVS to different background luminance. (*F2*) Weber-Fechner law [6, 9]. This law reveals important principle in psychophysics, which describes the logarithmic mapping between the magnitude of a physical stimulus and its perceived intensity. (*F3*) Texture masking [21]. Human eyes are more sensitive to perturbations on pixels in smooth regions than those in textured regions. (*F4*) Channel modulation [28, 31, 41]. Human vision has different sensitivity to the perturbations in different color channels.

- **Key observation.** These four factors clearly show that the perceptibility of human eyes is a **synthesized ability**; the perceptibility of human eyes is too complex to be reliably measured by a single metric.

Based on this key observation, we found that the limitations of the existing works are primarily due to the fact that they are all relying on a single metric. First, most of the existing works utilize distance metrics of $L_p$ norms ($L_0$, $L_2$ and $L_\infty$ norms) to measure imperceptibility. For example, Papernot et al. [24] generated adversarial examples based on a precise understanding of the mapping between inputs and outputs of DNNs by $L_2$ norm. The fast gradient sign method (FGSM) [34] computed one-step gradient to synthesize adversarial examples by $L_\infty$ norm. Carlini et al. [4] used three distance metrics $L_0$, $L_2$ and $L_\infty$ to quantify imperceptibility. Regarding the limitations of these works, it is obvious that they do not treat perceptibility of the human eyes as a synthesized ability. As a result, the measured perceptual loss could be conflicting with factors *F1-F4* for some images. Second, Luo et al. [21] defined perturbation sensitivity distance (PSD), taking the HVS into consideration, to measure perceptual loss between the benign examples and the adversarial ones. Zhao et al. [41] minimized perturbation size with respect to perceptual color distance in RGB space to generate large yet imperceptible adversarial perturbations. Unfortunately, these works again does not treat perceptibility of human eyes as a synthesized ability. As a result, the measured perceptual loss in [21] could be conflicting with factors *F1*, *F2*, and *F4*. And in [41], the measured perceptual loss could not consider factors *F1-F3*.

To effectively address the limitations of the existing methods, this work seeks to integrate JND, Weber-Fechner law, texture masking and channel modulation and design a **synthesized metric** to measure the perceptual loss caused by adversarial perturbations. To the best of our knowledge, this is the *first* metric ever proposed to measure perceptibility as a synthesized ability of human eyes in the context of generating adversarial examples. In order to design the synthesized metric, we first investigate how the four primary factors (JND, Weber-Fechner law, texture masking

and channel modulation) are complementary to each other. We then design a pixelwise multi-factor metric for measuring perceptual loss caused by perturbations.

In order to see whether the synthesized metric can let the generated adversarial examples "enjoy" much improved imperceptibility, the existing works could be classified into two categories: back-propagation [10, 20, 22, 23, 34] and forward-propagation [21, 24, 32]. Back-propagation utilizes output error and gradient information to compute adversarial perturbations, eg. L-BFGS [34], FGSM [10], etc. This type of attack results in whole pixel perturbations, thus, it is not suitable for optimizing pixelwise metrics. Forward-propagation computes a direct mapping from the input to the output to achieve an explicit adversarial goal. Representative studies include Jacobian-based saliency map attack (JSMA) [24], evolutionary algorithms [32], etc. Since JSMA utilizes the forward derivative to construct adversarial saliency maps, it can not be deployed in a black-box scenario. Evolutionary algorithms are a set of heuristic search methods, among which differential evolution has been proved to be a reliable tool to generate pixelwise perturbations [32]. However, if attackers fail to set the appropriate number of perturbed pixels, either the adversarial attack fails or the perturbations are unnecessarily large.

In order to overcome these limitations, assuming that the attacker does not have any inner information of the target DNNs, we apply a forward-propagation combining differential evolution and greedy approximation [19] to minimize the proposed pixelwise multi-factor metric. Greedy approximation allows GreedyFool to automatically find which pixels to perturb and what magnitude to modify effectively, guaranteeing successful black-box attacks with less perceptual loss.

In summary, this paper makes the following contributions:

- We investigate the HVS and design a synthesized metric to calculate the perceptual loss between the benign examples and the adversarial ones. The synthesized metric can let the generated adversarial examples "enjoy" much improved imperceptibility than the existing methods.

- To verify the advantages of the synthesized metric, We propose a pixelwise adversarial attack that is referred to as GreedyFool. GreedyFool attacks a target DNN only relying on its confidence. Thus, GreedyFool can attack more types of DNNs in a black-box manner. Since there are no budget constraints on the perceptual loss, GreedyFool achieves a 100% success rate for both non-targeted attacks and targeted attacks. The code is available on: https://github.com/LiuHwell/greedyfool.git.

The remainder of this paper is organized as follows. We

present the preliminaries in Sec. 2 and give the careful investigation into the HVS in Sec. 3. In Sec. 4, we give details about the design and the implementation of Greedy-Fool. The experimental evaluations are presented in Sec. 5. Finally, we conclude the paper in Sec. 6.

## 2. Preliminaries

### 2.1. Deep neural network

Deep neural networks [12, 14, 15, 30] are usually constructed with multiple neural network layers. A neural network layer consists of a set of perceptrons and each perceptron maps a set of inputs to output values with an activation function, e.g., Sigmoid, ReLU. DNNs can be formed in a chain.

$$f(x, \theta) = f^{(k)}(\dots f^{(2)}(f^{(1)}(x, \theta_1), \theta_2), \theta_k) \qquad (1)$$

where $f^{(i)}(x, \theta_i)$ is the function of the $i'$th layer of the network, where $i = 1, 2, ..., k$. $x$ is the input example and $\theta_i$ is the weight of the $i'$th layer.

Convolutional neural networks is one of the most widely used neural networks, which comprise convolutional layers, pooling layers and fully connected layers. In order to verify the performance of GreedyFool, we perturb the CIFAR-10 dataset over DenseNet, and the ImageNet dataset [27] over Inception V3 [33].

### 2.2. Adversarial attacks

Adversarial attacks [20, 22–24, 32] aim to mislead the DNN-based classifier to an incorrect label by adding small perturbations in the benign examples, even if the perturbations are barely recognizable by human eyes. Adversarial attacks would be formulated as a box-constrained optimization problem, that is,

$$\begin{aligned} \min \quad & \|\delta\|_p \\ s.t. \quad & f(X) = l \\ & f(X') = l' \\ & l \neq l' \\ & X' = X + \delta \in D \end{aligned} \qquad (2)$$

where a trained DNN model $f$ predicts the benign example $X$ and the adversarial example $X'$ into $l$ and $l'$. $\|\cdot\|$ denotes the perceptual loss between two examples, and $\delta$ denotes the perturbation matrix that is added to the benign example for synthesizing the adversarial example, which remains in the benign domain $D$.

### 2.3. $L_p$ norms

Definition of adversarial examples requires a distance metric to quantify similarity between the benign image and the adversarial image. Existing works mainly apply $L_p$ norms (e.g. $L_0$, $L_2$ and $L_\infty$) to measure the magnitude of the perturbation $\delta$ [4, 24, 34], that is,

$$\|\delta\|_p = \left( \sum_{i=1}^n |\delta_i|^p \right)^{\frac{1}{p}} \qquad (3)$$

where $L_0$ norm measures the number of pixels perturbed in an image. $L_2$ measures the Euclidean distance between the benign example and the adversarial example. $L_\infty$ norm denotes the maximum for all vector elements $|\delta_i|$ : $\|\delta\|_\infty = max(|\delta_i|)$. In these definitions, the small $L_p$ norm value indicates the greater imperceptibility of perturbations.

### 2.4. Differential evolution

Differential evolution [3,5,32] is a stochastic population-based algorithm for solving global optimization problems. The framework of differential evolution consists of the iteration of 3 operations: "one-to-one-spawning" selection, mutation, and crossover. The "one-to-one-spawning" selection mechanism replaces the parent solutions with the fitter child solutions. During each iteration, the algorithm mutates each candidate solution by mixing with other candidate solutions to create a trial candidate. There are several mutation strategies for creating trial candidates. Among these strategies, the "rand/1" strategy is proved to be feasible in one-pixel attack [32], which is defined in Eq. 4.

$$x_i(g+1) = x_{r1}(g) + F\left( x_{r2}(g) - x_{r3}(g) \right), r1 \neq r2 \neq r3 \qquad (4)$$

where $r1$, $r2$, and $r3$ are random numbers and $F$ is the scale parameter, which is set to 0.5. Three candidate solutions $x_{r1}(g)$, $x_{r2}(g)$ and $x_{r3}(g)$ from the $g'$th generation are randomly chosen to generate a new candidate solution $x_i(g+1)$ of the $g + 1'$th generation. Note that the crossover strategy is not included in our method as in [32].

## 3. Human visual system

The human visual system is a multichannel model with characteristics of multifrequency channel decomposition [11]. The sensitivity of human eyes to perturbations is affected by several factors. We investigate the HVS and summarize four primary factors that affecting human vision in the image domain: JND, Weber-Fechner law, texture masking and channel modulation.

### 3.1. Just noticeable distortion

Human eyes cannot perceive a stimulus below the just noticeable distortion [35, 38] threshold around a pixel in images. A larger value of the JND threshold indicates that more noise can be hidden. The visibility threshold of the

JND is formulated in Eq. 5, and shown as the curve in Fig. 1.

$$jnd(x,y) = \begin{cases} 17(1 - \sqrt{\frac{\bar{I}_{n \times n}(x,y)}{127}}) + 3, \text{if } \bar{I}_{n \times n}(x,y) \leq 127 \\ \frac{3}{128}\left(\bar{I}_{n \times n}(x,y) - 127\right) + 3, \text{otherwise} \end{cases}$$
$$(5)$$

where $\bar{I}_{n \times n}(x,y)$ denotes the maximum luminance among a $n \times n$ window at the coordinates $(x,y)$. Especially, we set $n = 3$ and calculate the maximum luminance of the $(x,y)$ pixel and its 8 neighbors.
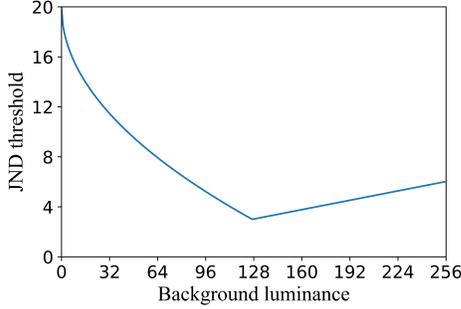


Figure 1. Illustration of the JND. When the maximum luminance is 127, the JND threshold reaches a minimum 3.

Figure 1 reveals the nonlinear relationship between the JND threshold and the maximum background luminance. Even if the magnitude of perturbations is the same, there are significant differences in the perceived intensity of human eyes in different luminance backgrounds.

### 3.2. Weber-Fechner law

Weber-Fechner law [6, 9] has been proposed in the field of psychophysics to quantify relationships between any stimulus and the perceived response by individuals. The Weber law asserts that the just noticeable stimulus difference $\Delta I$ maintains a constant ratio with respect to the intensity of the comparison stimulus $I$. On the assumption that the difference threshold represents a unit change in sensation $\Delta S$, Fechner defined Weber's law as:

$$\Delta S = k \frac{\Delta I}{I} \qquad (6)$$

Integrating this formula, a logarithmic relation called the psychophysical law is generated, as shown in Eq. 7.

$$S = k \ln I + C \qquad (7)$$

where $k$ and $C$ are hyperparameters, $I$ denotes the magnitude of a physical stimulus and $S$ denotes the corresponding perceived intensity. The Weber-Fechner law revealed the logarithmic mapping between the magnitude of a physical stimulus and its perceived intensity, rather than the absolute magnitude value as demonstrated in [21].



Figure 2. Perceptual loss of the perturbations with the same magnitude in different regions. The green line box marks perturbations in perturbed images. (a) benign example, (b) perturbed image at textured regions, (c) perturbed image at smooth regions.

### 3.3. Texture masking

According to texture masking theory [21], human eyes are more sensitive to perturbations on pixels in smooth regions than those in textured regions. We add 3×3 window perturbations with 100 magnitudes in the smooth region and the textured region for comparison of the perceptual loss. Figure 2 plots the benign example, the perturbed image at textured regions and perturbed image at smooth regions from left to right. When perturbations with the same magnitude are added to the image, the perturbations in smooth regions are easily detected by human eyes, while those in textured regions are difficult to recognize.

The standard deviation is a commonly employed quantity for measuring the texture masking of an image, which is proven to be effective in evaluating the perceptual loss of adversarial examples. The paper computes the standard deviation of a pixel $p_i$ in a $n \times n$ window as shown in Eq. 8.

$$SD\left(p_i\right) = \sqrt{\frac{\sum_{p_i \in S_i}\left(p_i - \mu\right)^2}{n^2}} \qquad (8)$$

where $S_i$ is the set of pixels in the $n \times n$ window, and $\mu$ comprise the average values of pixels in the region. The standard deviation $SD(p_i)$ is the variance of pixel $p_i$ among $S_i$. In this paper, we set $n = 3$ and calculate the standard deviation of pixel $p_i$ and its 8 neighbors.

### 3.4. Channel modulation

According to spectral sensitivity theory [28, 31], human eyes consist of red, green and blue cone cells, and the cone is sensitized to different ranges of wavelengths to provide a range of color perception. The human eyes are the most sensitive to green, followed by red, and the least sensitive to blue. We conduct perturbations with the same magnitude in three color channels. As shown in Fig 3, the test proves that there are differences in the sensitivity of human eyes to color channels.

We introduce channel modulation to quantify the weight of perturbations in three color channels. Channel modulation refers to a constrained linear combination of red,
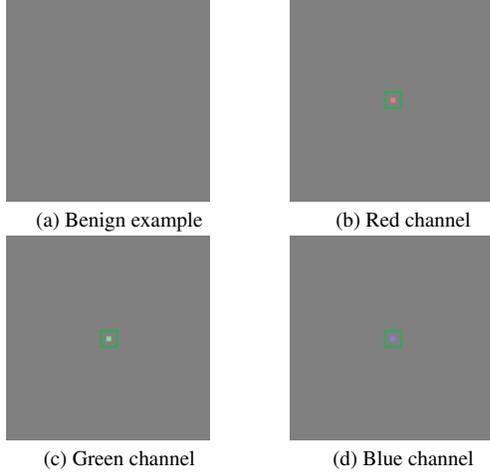
Figure 3. Perceptual loss of the perturbations with the same magnitude in different channels. The green line box marks perturbations in perturbed examples.

green and blue channels based on decolorization theory [31], which is formulated in Eq. 9.

$$
\begin{aligned}
\varpi &= \lambda_r I_r + \lambda_g I_g + \lambda_b I_b \\
s.t. \quad & \lambda_r + \lambda_g + \lambda_b = 1 \\
& \lambda_r \geq 0, \lambda_g \geq 0, \lambda_b \geq 0
\end{aligned}
\tag{9}
$$

where $I_r$, $I_g$ and $I_b$ are the red channel, green channel and blue channel respectively, and $\varpi$ is the result of the channel modulation. The non-negative numbers $\lambda_r$, $\lambda_g$ and $\lambda_b$ are channel weights that sum to 1. In the classical RGB2GRAY conversion model [31], the weights are fixed as $\lambda_r = 0.299, \lambda_g = 0.587, \lambda_b = 0.114$.

## 4. Methodology

Adversarial attacks could be formulated as an optimization problem with constraints. This problem involves the definition of the perceptual loss and the solution for the optimal problem. In this section, we describe a definition of the multi-factor perceptual loss and present a forward-propagation that combines differential evolution with greedy approximation

### 4.1. Multi-factor perceptual loss

In the investigation into the HVS, we discover that the perceptual loss is primarily affected by four factors: JND, Weber-Fechner law, texture masking and channel modulation. In computer vision, the application of Weber-Fechner law needs to consider the difference in the perceived intensity of the human eyes with different luminance backgrounds. The visibility of the human eyes to perturbations of images has a positive correlation with a stimulus and a

negative correlation with JND. Thus, we combine Eqs. 5 and 6 to redefine this correlation, as shown in Eq. 10.

$$
\Delta Ps = k \frac{\Delta I}{JND(I)}
\tag{10}
$$

where $k$ is a constant according to Weber-Fechner law, $\Delta I$ is the just noticeable stimulus difference and $\Delta Ps$ is the corresponding change in visual sensation. Since $JND(I)$ is a discrete piecewise function, we sum over Eq. 10 rather than integrate, as shown in Eq. 11.

$$
Ps(I) = k \sum_{i=I_0}^{I-1} \frac{1}{JND(i)} + C
\tag{11}
$$

where $Ps(I)$ is the perceptual stimulus, $I_0$ is the background luminance, and $k$ and $C$ are constants.

To construct the mapping between a physical stimulus and its perceptual stimulus, we define the perceptual stimulus in the interval [0, 255]. Obviously, a perturbation of magnitude 0 does not cause any perceptual perception to the human eyes and a perturbation of magnitude 255 could cause the most noticeable perceptual perception to the human eyes. Thus, the following equation could be constructed to compute the constants $k$ and $C$, that is,

$$
\begin{cases}
Ps(0) = C = 0, \text{if } I = 0 \\
Ps(255) = k \sum_{i=0}^{255} \frac{1}{JND(i)} + C = 255, \text{if } I = 255
\end{cases}
\tag{12}
$$

By solving Eq. 12, we obtain the parameters $C = 0$ and $k = \frac{256}{\sum_{i=0}^{255} \frac{1}{JND(i)}}$.

Texture masking reflects that perturbations in the region with high standard deviation are more imperceptible than those in the region with low standard deviation. The visibility of human eyes to perturbations has a negative correlation with the standard deviation of the region in an image. Furthermore, considering channel modulation, we propose a synthesized metric for evaluating the perceptual loss as follows.

$$
IntegLoss(p_i) = \sum_{c \in (r,g,b)} \lambda_c \frac{Ps(p_i)}{SD_c(p_i)}
\tag{13}
$$

Adversarial attacks usually add multiple pixel perturbations for the success rate. The perceptual loss between the benign image and the adversarial image is the sum of all pixelwise perceptual losses. Therefore, we sum all pixelwise perceptual losses as follows.

$$
MulFactorLoss(X, X') = \sum_{i=1}^{N} \sum_{c \in (r,g,b)} \lambda_c \frac{Ps_c(p_i)}{SD_c(p_i)}
\tag{14}
$$

where $N$ is the number of perturbed pixels in a benign image. $MulFactorLoss(X, X')$ denotes the perceptual loss

that integrates *JND*, Weber-Fechner law, texture masking and channel modulation between the benign example $X$ and the adversarial example $X'$. The small $MulFactorLoss$ value indicates the high perceptual similarity between the benign example and the adversarial example.

### 4.2. Pixelwise objective function

State-of-the-art adversarial attacks should allow DNNs to give the wrong output with a high confidence score by adding as few perturbations as possible. Therefore, we should choose pixels that can reduce the confidence of DNNs in the true class or increase that in the target class with the less perceptual loss. The pixelwise objective function, which is referred to as the perturbation priority, is defined to estimate the effect of perturbing a pixel as follows.

$$PertPriority(p_i) = \zeta \frac{P_t(X) - P_t(X')}{MulFactorLoss(p_i)} \quad (15)$$

where $\zeta$ is a control parameter. $P_t$ denotes the probability that the example belongs to the label $t$. The adversarial example $X'$ is synthesized by changing a pixel $p_i$ of benign example $X$. When $t = l$, $\zeta = 1$ and non-targeted attacks are executed, otherwise $\zeta = -1$ and targeted attacks are executed. The perturbation priority quantifies the effects of the current pixel $p_i$ perturbation on the confidence of the DNN-based classifier in the target class.

As can be seen from Eq. 15, this objective function only relies on the DNN's confidence on the test image. Thus, GreedyFool can attack more types of DNNs in a black-box manner.

### 4.3. Implementation

To determine the pixels with high priority, the adversary has to choose which pixels to modify and what magnitudes to add. We encode the pixelwise perturbation into an array as a candidate solution $(x, y, r, g, b)$, which contains five elements: $x - y$ coordinates and RGB value of the perturbation. A brute-force approach has to search all dimensions and pixel values for the optimal value, which could take a prohibitively long time.

To reduce the search time, we introduce differential evolution to solve the optimal pixelwise objective function. GreedyFool sets the total population size to 200 and the number of generations to 60, resulting in 12,000 candidate solutions with perturbation priority. Greedy approximation is utilized to automatically obtain a set of perturbed pixels and synthesize the imperceptible adversarial example. The implementation of GreedyFool is presented as follows.

Step 1. Encoding the perturbation into candidate solution $(x, y, r, g, b)$, and randomly initializing 200 candidate solutions;

Step 2. Executing differential evolution during 60 generations to calculate the priority of candidate solutions;

| Method | GreedyFool | [34] | [10] | [1] |
|---|---|---|---|---|
| *MulFactorLoss* | 49.81 | 432.16 | 365.87 | 137.35 |

Table 1. *MulFactorLoss* values of adversarial examples on imageNet

Step 3. Greedy approximation is executed to choose a set of candidate solutions with the highest priority as adversarial perturbations. These perturbations are added to the benign image to synthesize its adversarial example.

## 5. Experimental evaluation

### 5.1. Experimental setup

In this experiment, we implement GreedyFool in Python and conduct the adversarial attack in ImageNet over Inception V3, and CIFAR-10 over DenseNet. Due to the size 224×224 of the preprocessed ImageNet and 32×32 of CIFAR-10, and the pixel value interval [0, 255], the population of the differential evolution is initialized by uniform distributions $U(1, 224)$ for ImageNet and $U(1, 32)$ for CIFAR-10 to generate the $x - y$ coordinates and Gaussian distributions $N(\mu = 128, \sigma = 127)$ for pixel RGB values. The fitness function of the differential evolution is set to a pixelwise objective function in Eq. 15. GreedyFool runs on an Intel Core I5 CPU 2.30 GHz, NVIDIA GeForce GTX 1060 and 8.0 GB of RAM computer that run on Windows 10 and Spyder (Python 3.6).

### 5.2. Adversarial examples on ImageNet

We evaluate GreedyFool and compare it with 3 state-of-the-art adversarial attacks (L-BFGS [34], FGSM [10], Color-and-Edge-Aware [1]) on the Inception V3 model for ImageNet classification task. Figure 4 shows the comparison of adversarial examples. Close inspection reveals perturbation artefacts in the sky near the water's surface for all methods except GreedyFool. Table 1 lists *MulFactorLoss* values of adversarial examples in Fig. 4. The adversarial example generated by GreedyFool presents less visibility and *MulFactorLoss* value than those generated by the other 3 attacks.

### 5.3. Human eye evaluation on CIFAR-10

To clearly present the details of the pixelwise perturbations, we generated adversarial examples on the CIFAR-10 dataset and carry out human eye tests. Similar to GreedyFool, both One-pixel attack [32] and PSD attack [21] belong to black-box pixelwise attack that use forward-propagation. Therefore, all three attacks are carried out on the same, random set of images to compare the imperceptibility of adversarial perturbations. We present 10 classes of images in
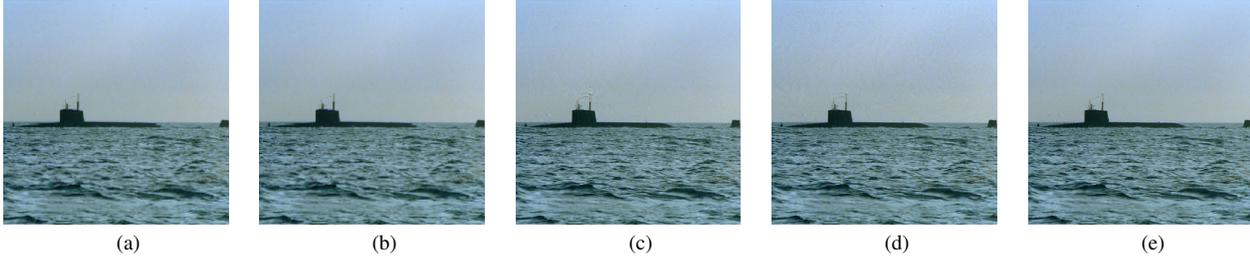
Figure 4. Comparison of adversarial examples generated by GreedyFool, L-BFGS, FGSM, and Color-and-Edge-Aware on the same target label breakwater.

Fig. 5, which include benign images and their adversarial images against DenseNet.

Imperceptibility is a subjective feeling that varies among individuals. To evaluate the visibility of adversarial perturbations, we design a scoring system based on the human eye evaluation and conduct an extensive user study. In this test, we generate 150 images for each attack over DenseNet, which are from 15 benign images targeting for 10 classes. We recruit 60 participants to score the similarity between benign images and their adversarial examples generated randomly by one of the aforementioned attacks. For each trial, the benign images are shown on a screen at a fixed size and the order of adversarial images is random. The participants are required to give a score from 0 to 10 in 5 seconds. A higher score denotes more noticeable adversarial perturbations. A score of 0 means no difference, while a score of 10 means a complete difference. We calculate the average score for each class of images from the human eye evaluation and plot them in Fig. 6 for comparison.

As shown in Fig. 6, GreedyFool has the lowest score in any classes. One-pixel attack has the highest score in most classes, because human eyes can easily detect perturbed pixels. The results show that adversarial examples synthesized by GreedyFool have higher perceptual similarity to their benign images. It indicates that GreedyFool generates more imperceptible adversarial examples than other pixelwise methods.

Among the three types of attacks, One-pixel attack hardly considers the HVS. PSD attack considers the intensity and texture masking of benign images. GreedyFool combines more abundant and effective factors to construct an objective function of adversarial attacks, which fully considers the HVS. The human eye experiment demonstrates that the imperceptibility of adversarial examples could be further improved by combining multiple HVS metrics.

## 5.4. Perceptual loss results

A good perceptual loss should reliably reflect the perceptual similarity between the benign example and the adversarial example. Table 2 lists *MulFactorLoss* results of the adversarial examples generated by these three methods. The results show that GreedyFool could generate adversarial examples with smaller *MulFactorLoss* values than PSD attack and One-pixel attack, which are consistent with the scores obtained from the human eye evaluation.

| Method | [32] | [21] | GreedyFool |
|---|---|---|---|
| airplane | 2.98 | 1.32 | 0.34 |
| automobile | 0.01 | 24.64 | 0.01 |
| bird | 11.09 | 3.45 | 2.41 |
| cat | 7.88 | 2.48 | 0.65 |
| deer | 7.32 | 0.67 | 0.40 |
| dog | 4.04 | 0.53 | 1.01 |
| frog | 10.87 | 2.84 | 1.55 |
| horse | 7.06 | 0.28 | 0.15 |
| ship | 13.98 | 0.17 | 0.04 |
| truck | 7.19 | 0.17 | 0.06 |

Table 2. Multi-factor perceptual loss results of adversarial examples over DenseNet

Table 3 lists the perceptual loss results of the adversarial examples generated by GreedyFool in Fig. 5 The results show that $L_p$ norms cannot always reliably quantify the imperceptibility of adversarial perturbations. The multi-factor perceptual loss *MulFactorLoss* is a more reliable metric than $L_p$ norms to reflect the perceptual similarity between the benign example and the adversarial example.

## 5.5. Misclassification ratio

According to the purpose of the tasks, adversarial attacks could be classified as non-targeted attacks and targeted attacks. Non-targeted attacks aim to mislead the DNN to classify images into any wrong classes, while targeted attacks aim to mislead the DNN to classify images into a specific target class.

We test 1,000 random images to compute the misclassification ratio for both non-targeted and targeted attacks. The results demonstrate that GreedyFool achieves 100% success
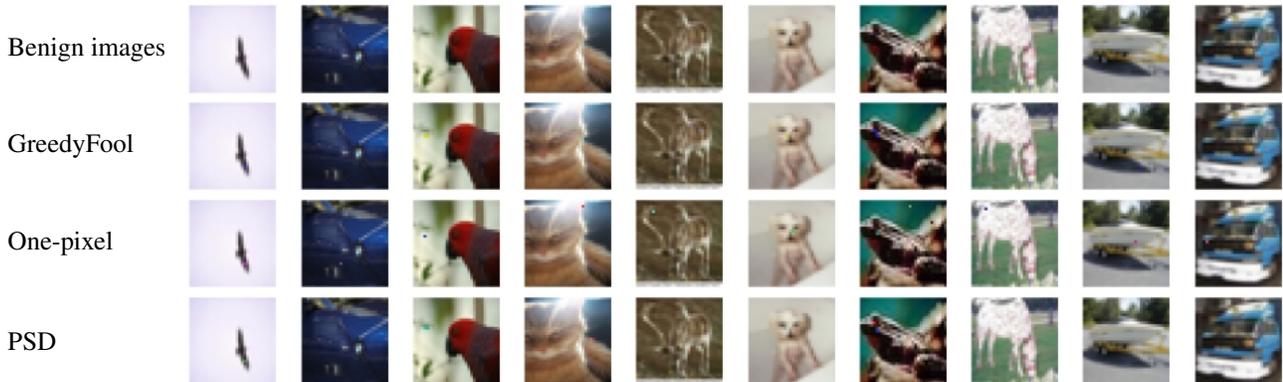
Figure 5. Adversarial images synthesized by different attack methods against DenseNet in CIFAR-10. Adversarial examples in the second row crafted by our method are much more imperceptible than other examples from the following rows.
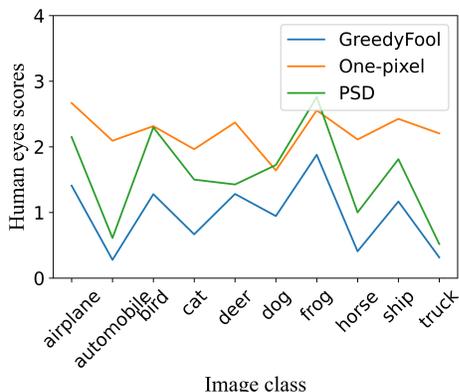


Figure 6. Average scores comparison from human eyes evaluation between our attacks with existing attacks on DenseNet.

| Perceptual loss | $L_0$ | $L_2$ | $L_\infty$ | $MulFactorLoss$ |
| --- | --- | --- | --- | --- |
| airplane | 1 | 210 | 210 | 0.34 |
| automobile | 1 | 2 | 2 | 0.01 |
| bird | 2 | 397 | 203 | 2.40 |
| cat | 1 | 125 | 125 | 0.65 |
| deer | 1 | 69 | 69 | 0.39 |
| dog | 1 | 93 | 72 | 1.01 |
| frog | 2 | 462 | 247 | 1.55 |
| horse | 1 | 75 | 75 | 0.15 |
| ship | 1 | 5 | 5 | 0.04 |
| truck | 1 | 47 | 47 | 0.06 |

Table 3. Perceptual loss results of adversarial examples generated by GreedyFool over DenseNet

rate for both non-targeted attacks and targeted attacks.

When the perturbations are sufficient, any images could be misclassified as a specified target class by DNNs. However, adversarial attacks do not always achieve a 100% suc-

cess rate because some of them need to establish the box-constrained parameters to limit perturbations. GreedyFool utilizes greedy approximation to release this constraint, which could automatically generate imperceptible adversarial examples in a 100% success rate.

## 5.6. Computation cost

Computation cost refers to the running time for attackers to synthesize an adversarial example, which is employed to evaluate the attack time cost. The computation cost of GreedyFool is affected by several factors, including the feedback time of the DNN model, convergence speed of differential evolution, number of greedy approximations, running environment, etc.

The experiment chooses 30 random images from CIFAR-10 to test the actual running time of GreedyFool. On average, GreedyFool consumes approximately 38.44 seconds to carry out an adversarial attack over DenseNet on CIFAR-10, and 94.62 seconds over Inception V3 on ImageNet.

## 6. Conclusion

Adversarial attack against neural networks is a serious threat to safety-critical systems. Existing works lack sufficient consideration of the HVS, which produces noticeable artifacts in adversarial examples. In order to obtain sufficient imperceptibility, we investigate the HVS and identify four primary factors affecting the perceptibility of the human eyes: JND, Weber-Fechner law, texture masking and channel modulation.

Based on these factors, we design a pixelwise multi-factor metric to define the perceptual loss between benign images and adversarial examples. To test the multi-factor metric, we propose a black-box approach that is referred to as GreedyFool to generate adversarial examples using forward-propagation. We implement GreedyFool and con-

duct the adversarial attack in ImageNet on the Inception V3, CIFAR-10 on DenseNet. The experimental results show that GreedyFool has greater imperceptibility than state-of-the-art pixelwise methods, which achieves a 100% success rate in a black-box manner.

# References

[1] Robert Bassett and Mitchell Graves. Color and edge-aware adversarial image perturbations. *CoRR*, abs/2008.12454, 2020. 6

[2] Anand Bhattad, Min Jin Chong, Kaizhao Liang, Bo Li, and David A. Forsyth. Unrestricted adversarial examples via semantic manipulation. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 1

[3] Fabio Caraffini, Anna V. Kononova, and David Corne. Infeasibility and structural bias in differential evolution. *Inf. Sci.*, 496:161–179, 2019. 3

[4] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 39–57. IEEE Computer Society, 2017. 1, 2, 3

[5] Pinar Civicioglu and Erkan Besdok. A conceptual comparison of the cuckoo-search, particle swarm optimization, differential evolution and artificial bee colony algorithms. *Artificial intelligence review*, 39(4):315–346, 2013. 3

[6] Stanislas Dehaene. The neural basis of the weber-fechner law: a logarithmic mental number line. *Trends in Cognitive Sciences*, 7(4):145–147, 2003. 2, 4

[7] Xingping Dong, Jianbing Shen, Dongming Wu, Kan Guo, Xiaogang Jin, and Fatih Porikli. Quadruplet network with one-shot learning for fast visual object tracking. *IEEE Trans. Image Process.*, 28(7):3516–3527, 2019. 1

[8] Zehao Dou, Stanley J. Osher, and Bao Wang. Mathematical analysis of adversarial attacks. *CoRR*, abs/1811.06492, 2018. 1

[9] Jan Drösler. An n-dimensional weber law and the corresponding fechner law. *Journal of Mathematical Psychology*, 44(2):330–335, 2000. 2, 4

[10] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 1, 2, 6

[11] Yu S Gulina and V Ya Kolyuchkin. Experimental investigations of a model of the human visual system. *Optics and Spectroscopy*, 127(4):675–683, 2019. 3

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016. 1, 3

[13] Shota Horiguchi, Sosuke Amano, Makoto Ogawa, and Kiyoharu Aizawa. Personalized classifier for food image recognition. *IEEE Trans. Multim.*, 20(10):2836–2848, 2018. 1

[14] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 7132–7141. Computer Vision Foundation / IEEE Computer Society, 2018. 1, 3

[15] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2261–2269. IEEE Computer Society, 2017. 1, 3

[16] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 125–136, 2019. 1

[17] Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense against unseen threat models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 1

[18] Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. Textbugger: Generating adversarial text against real-world applications. In *26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019*. The Internet Society, 2019. 1

[19] Hui Liu, Bo Zhao, Linquan Huang, Jiabao Guo, and Yifan Liu. Foolchecker: A platform to evaluate the robustness of images against adversarial attacks. *Neurocomputing*, 412:216–225, 2020. 2

[20] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. 1, 2, 3

[21] Bo Luo, Yannan Liu, Lingxiao Wei, and Qiang Xu. Towards imperceptible and robust adversarial example attacks against neural networks. In Sheila A. McIlraith and Kilian Q. Weinberger, editors, *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 1652–1659. AAAI Press, 2018. 1, 2, 4, 6, 7

[22] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 86–94. IEEE Computer Society, 2017. 1, 2, 3

[23] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method

to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2574–2582. IEEE Computer Society, 2016. 1, 2, 3

[24] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *IEEE European Symposium on Security and Privacy, EuroS&P 2016, Saarbrücken, Germany, March 21-24, 2016*, pages 372–387. IEEE, 2016. 1, 2, 3

[25] Haonan Qiu, Chaowei Xiao, Lei Yang, Xinchen Yan, Honglak Lee, and Bo Li. Semanticadv: Generating adversarial examples via attribute-conditional image editing. *CoRR*, abs/1906.07927, 2019. 1

[26] Erwin Quiring, Daniel Arp, and Konrad Rieck. Forgotten siblings: Unifying attacks on machine learning and digital watermarking. In *2018 IEEE European Symposium on Security and Privacy, EuroS&P 2018, London, United Kingdom, April 24-26, 2018*, pages 488–502. IEEE, 2018. 1

[27] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.*, 115(3):211–252, 2015. 3

[28] J Schnapf, Timothy Kraft, and Denis Baylor. Spectral sensitivity of human cone photoreceptors. *Nature*, 325:439–41, 01 1987. 2, 4

[29] Ali Shafahi, W. Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. 1

[30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 1, 3

[31] Yibing Song, Linchao Bao, Xiaobin Xu, and Qingxiong Yang. Decolorization: is rgb2gray() out? In Baoquan Chen and Andrei Sharf, editors, *SIGGRAPH Asia 2013, Hong Kong, China, November 19-22, 2013, Technical Briefs*, pages 15:1–15:4. ACM, 2013. 2, 4, 5

[32] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Trans. Evol. Comput.*, 23(5):828–841, 2019. 1, 2, 3, 6, 7

[33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer Society, 2016. 3

[34] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada,*

*April 14-16, 2014, Conference Track Proceedings*, 2014. 1, 2, 3, 6

[35] Zhibo Wang, Mengkai Song, Siyan Zheng, Zhifei Zhang, Yang Song, and Qian Wang. Invisible adversarial attack against deep neural networks: An adaptive penalization approach. *IEEE Trans. Dependable Secur. Comput.*, 18(3):1474–1488, 2021. 1, 2, 3

[36] Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of descriptor and generator networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(1):27–45, 2020. 1

[37] Mingfu Xue, Chengxiang Yuan, Can He, Jian Wang, and Weiqiang Liu. Naturalae: Natural and robust physical adversarial examples for object detectors. *J. Inf. Secur. Appl.*, 57:102694, 2021. 1

[38] Xiaokang Yang, W. S. Ling, Zhongkang Lu, Ee Ping Ong, and Susu Yao. Just noticeable distortion model and its applications in video coding. *Signal Process. Image Commun.*, 20(7):662–680, 2005. 2, 3

[39] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing [review article]. *IEEE Comput. Intell. Mag.*, 13(3):55–75, 2018. 1

[40] Xinze Zhang, Junzhe Zhang, Zhenhua Chen, and Kun He. Crafting adversarial examples for neural machine translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 1967–1977. Association for Computational Linguistics, 2021. 1

[41] Zhengyu Zhao, Zhuoran Liu, and Martha A. Larson. Towards large yet imperceptible adversarial image perturbations with perceptual color distance. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 1036–1045. Computer Vision Foundation / IEEE, 2020. 1, 2