# Feature-Filter: Detecting Adversarial Examples through Filtering out Recessive Features

Hui Liu   Bo Zhao   Minzhi Ji   Yuefeng Peng   Jiabao Guo
School of Cyber Science and Engineering, Wuhan University
Wuhan, 430072 China

{liuh824, zhaobo, jiminzhi, yuefengpeng, garbo_guo}@whu.edu.cn

Peng Liu
College of Information Sciences and Technology Pennsylvania State University
PA, 16801 US

pliu@ist.psu.edu

## Abstract

*Deep neural networks (DNNs) are under threat from adversarial example attacks. The adversary can easily change the outputs of DNNs by adding small well-designed perturbations to inputs. Adversarial example detection is a fundamental work for robust DNNs-based service. Adversarial examples show the difference between humans and DNNs in image recognition. From a human-centric perspective, image features could be divided into dominant features comprehensible to humans, and recessive features incomprehensible to humans, yet exploited by DNNs. In this paper, we reveal that imperceptible adversarial examples are the product of recessive features misleading neural networks, and the adversarial attack enriches these recessive features. The imperceptibility of the adversarial examples indicates that the perturbations enrich recessive features, yet hardly affect dominant features. Therefore, adversarial examples are sensitive to filtering out recessive features, while benign examples are immune to such operation. Inspired by this idea, we propose a label-only adversarial detection approach that is referred to as feature-filter. Feature-filter utilizes discrete cosine transform (DCT) to approximately separate recessive features from dominant features and gets a filtered image. A comprehensive user study demonstrates that the DCT-based filter can reliably filter out recessive features from the test image. By only comparing DNN's prediction labels on the input and its filtered version, feature-filter can real-time detect imperceptible adversarial examples at high accuracy and few false positives.*

## 1. Introduction

Deep neural networks (DNNs) [15–17, 37] achieve state-of-the-art performance on many artificial intelligence tasks, including safety-critical systems like facial biometric systems [3], intrusion detection system [26], etc. However, a large number of studies have shown that attackers can easily change the outputs of DNNs by the inputs with well-designed perturbations (adversarial examples) [41]. The malicious action is called "adversarial attack" [2, 14, 21, 22, 27, 28, 39, 48].

Adversaries could modify any pixel value in the image to generate adversarial examples. In the physical world, however, they don't want modifications that can easily irritate the human eye. For example, autonomous vehicles deploy extensive neural networks to construct their perceptual systems [7, 45]. Adversaries could try to affect the decision-making of perceptual systems by using adversarial traffic signs [23]. If these modifications in traffic signs can be easily detected by the human eye, the traffic police would replace them in time and the adversarial attack would fail before it had even begun. Therefore, in order to achieve this goal, adversaries must ensure that these modifications are imperceptible to the human eye. This constraint of imperceptibility is widespread in real-world scenarios.

The phenomenon of adversarial examples has attracted a great number of interests in the research community in recent years. To hardening DNN systems, A wide range of proactive defense approaches [20, 30, 32, 36, 43] have been proposed in image classification, but those defenses can be evaded by emerging attack approaches [6]. If there are no known intrinsic properties that distinguish adversarial examples from benign examples, these proactive adversarial defenses are extremely challenging [8].

As a fundamental work for robust DNNs, detecting adversarial examples may be as important as adversarial defenses. Adversarial detection aims to distinguish adversarial examples from benign ones by intrinsic properties. To discover these intrinsic properties, let us analyze the reasons for the existence of imperceptible adversarial examples in the image recognition.

There has been a debate on why the adversarial examples exist. Previous work has proposed several explanations for the phenomenon of adversarial examples, ranging from finite-sample overfitting to high-dimensional statistical properties [4, 25, 35]. However, Andrew et al. [18] argued these theories were often unable to fully capture behaviors, and first categorized image features from a human-centric perspective. Inspired by this study, we observe that adversarial examples show the difference between humans and DNNs in the utilization of input features. For accurate classification, DNNs make full use of the input features, of which spaces are often unnecessarily large for humans. The human eyes recognize images only according to those features that are comprehensible to humans. As a human-centric phenomenon, the input features could be divided into two categories: dominant features comprehensible to humans, and recessive features incomprehensible to humans. When adversarial perturbations enough enrich recessive features and hardly lose dominant features, they would mislead the output of DNNs and not result in noticeable artifacts to the human eyes. Thus, we claim that the existence of imperceptible adversarial examples is due to the neural network's use of these recessive features.

Since adversarial examples generated by state-of-the-art attack approaches have rich recessive features to mislead DNN models, they would be much more sensitive to filter out these recessive features than benign examples. A natural detection strategy of adversarial examples is driven by comparing the model's predictions on the original input and its filtered version, as depicted in Fig. 1. Where, the filter filters out recessive features of the original image as a filtered image. If the original image and its filtered version are classified by the DNN model into different categories, this image would be judged as an adversarial example.

A challenging task is how to reliably filter out the recessive features from the original image. A key idea is inspired by classical JPEG compression based on discrete cosine transform (DCT) [11, 12, 46]. As depicted in Fig. 2, the compressed image is still clearly recognizable to humans, even though the feature space is only 1/10 the size of the original image. It indicates that a large number of dominant features are retained, and the contents filtered out almost exclusively include recessive features. Therefore, DCT could be employed as the filter in Fig. 1 to filter out recessive features from input features.
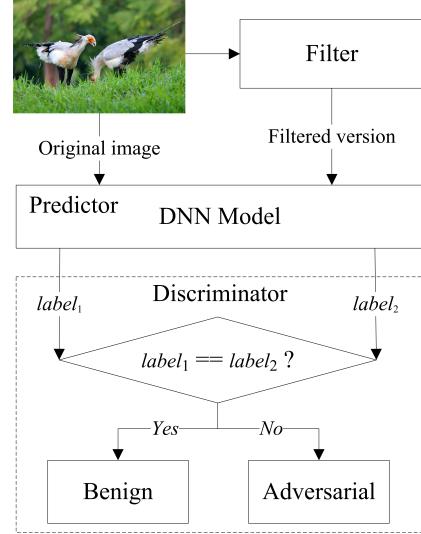
In summary, this paper makes the following contributions:



Figure 1. Feature-filter framework for adversarial detection.
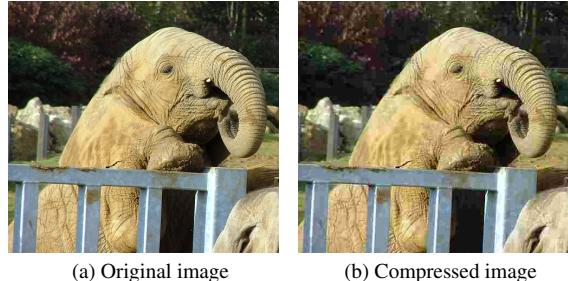


(a) Original image      (b) Compressed image

Figure 2. JPEG compression.

- We analyze the reason for the existence of imperceptible adversarial examples, and reveal that the imperceptible adversarial example is the product of abundant recessive features misleading neural networks.

- To clarify the above explanation, we propose a label-only adversarial detection approach that is referred as feature-filter. Feature-filter can identify adversarial examples by only comparing prediction labels of the DNN on the original image and its filtered version.

The remainder of this paper is organized as follows. Section 2 introduces the related work about adversarial attacks, defenses and detections. Section 3 is the preliminaries about the proposed detection. In Sec. 4, we give details about the design and the implementation of the feature-filter. Section 5 evaluates the performance of the DCT-based filter and our detector. Finally, we conclude the paper in Sec. 6.

## 2. Related work

In this section, we briefly review existing works on adversarial attacks, defenses and detections.

**Adversarial attacks:** The adversarial example is first discovered in exploring properties of neural networks. Szegedy et al. [41] found that the DNN would misclassify a well-designed image, in which certain perturbations were hardly perceptible to the human eyes. They designed a box-constrained L-BFGS to compute an approximate minimum perturbation matrix.

Many improved works have been subsequently proposed by modifying the adversarial objective function. For example, the fast gradient sign method (FGSM) [14] updates each pixel along the direction of the sign of gradient and performs a fast one-step attack. The fast gradient value method (FGVM) [33] replaces the sign of the gradient with the raw gradient, and generates adversarial examples with a larger local difference. Despite their efficiency, these methods provides only a coarse approximation of the optimal perturbations, resulting in unnecessarily large artifacts.

To improve the imperceptibility of perturbations, Deep-Fool [28] designed an accurate method by moving the input image across its closest decision boundary. Adversarial examples generated by DeepFool have less perturbations in the input image than those generated by L-BFGS and FGSM. Inspired by DeepFool, an universal adversarial attack [27] was proposed to reveal important geometric correlations among the high-dimensional decision boundary of DNNs, which proved the transferability of adversarial examples across neural networks.

Among existing adversarial attacks, C&W [6] attack performs high resilience against the most defense and detection methods. Especially, $L_2$ attack among them is argued to be the most effective to date in white-box threat model. Therefore, most adversarial detectors are evaluated by testing in adversarial examples generated by C&W attack.

**Adversarial defenses:** A wide variety of proactive defense approaches have been proposed to harden the neural network. Based on a comprehensive summary, adversarial defenses are grouped into two broad categories: adversarial training and gradient masking [29].

Adversarial training extends the training dataset by adversarial examples with the ground-truth labels and retrains the neural network using this extended dataset. Goodfellow et al. [14] claimed that adversarial training could provide an additional regularization benefit and increase the robustness of neural networks. Although this type of defense works against attacks on which the model is trained, it faces potential risks for new attacks.

Gradient masking is an effective defensive strategy against gradient-based adversarial attacks. These attack methods compute first-order derivatives to estimate the sensitivity of the model with respect to inputs. Papernot et al. [31] introduced "defensive distillation" to defend DNNs against adversarial examples. To hide the model's gradients, network distillation replaces its softmax layer with a "softer" softmax function. In fact, gradients are not essential knowledge for attackers, who easily evade gradient masking using a transferability-based black-box attack [27].

**Adversarial detections:** Due to the difficulty of adversarial defenses, recent works focus on reactive adversarial detections, which distinguish adversarial images from benign images. Adversarial detections could be grouped into three categories: statistical measurement [9, 13, 47], secondary classifier [24, 38] and input transformation [1, 19, 42, 44].

Statistical measurement aims to find the statistical differences in the high-dimensional features extracted from DNN layers. NNIF [9] utilized influence functions and nearest neighbors to construct k-NN features, which are employed to train a logistic regression for adversarial detection. ML-LOO [47] computes the contribution of features in the prediction of model to construct feature attribution map, and then measures 3 metrics (standard deviation, median absolute deviation and interquartile range) of statistical dispersion in feature attribution map. The experiments revealed a significant difference in the distributions of 3 metrics between benign and adversarial images. Due to the intrinsically unperceptive nature of adversarial examples, statistical measurement seems unlikely to be effective to detect adversarial examples.

Adversarial detections is essentially binary classification problem. A nature idea is to train secondary classifiers to detect if the input is benign or adversarial. NIC [24] trains a model for each layer in the DNN to describe the distributions of the provenance and the activation value. By combining the results of all models, NIC makes a joint prediction to determine whether the input image is adversarial or not. In fact, the secondary classifiers also could be fooled, thus, they have the potential vulnerability to second-round adversarial attacks.

The basic idea of input transformation is to measure the disagreement of the target model in predicting the input and its transformed versions. Xu et al. [44] proposed a detection strategy, feature squeezing, to squeeze feature space available to attackers. However, they did not specify which features should be squeezed. Kantaros et al. [19] observed that adversarial examples are sensitive to lossy compression transformations. They proposed VisionGuard to real-time detect large-scale adversarial examples. VisionGuard measured similarity of softmax outputs between the test image and its compressed version using the K-L divergence measure. When the the K-L divergence is greater than a threshold, the test image is considered an adversarial input. VisionGuard determined the detection threshold based on the training dataset and its adversarial dataset, which makes

it inapplicable to data with different distributions. Tian et al. [42] exploited a set of rotation operations to yield several transformed versions and analyzed the classification results of these transformed images to determine if the test image is adversarial. Since input transformation can not change neural networks, it is inexpensive and complementary to other defenses and detections.

## 3. Preliminaries

### 3.1. Neural network

A neural network [15–17,37,40] has a hierarchical structure, in which each layer is composed of a set of perceptrons. Perceptrons map inputs to output values with a non-linear activation function, e.g. sigmoid, tanh, and ReLU. The function of a neural network is formed in a chain.

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

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

Convolutional neural network is one of the most widely used neural networks, which consists of alternating convolutional layers and pooling layers. Convolutional neural network has performed incredible successes in image recognition. As a representative network, Inception V3 [40] has 23,851,784 parameters and performs 93.70% top-5 accuracy for the ImageNet datase.

### 3.2. Adversarial attack

Adversarial attack [2,14,21,22,27,28,39,48] is a unique form of attack in deep learning, in which attackers are able to change DNNs' outputs by adding imperceptible perturbations to inputs. Adversarial attack would be formulated as an optimization problem with imperceptibility as the main constraint, that is,

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

where $\delta$ denotes the perturbation matrix that is added to the benign example $x$ for synthesizing the adversarial example $x'$, which remains in the benign domain $D$, and a trained DNN model $f$ predicts $x$ and $x'$ into different labels ($l \neq l'$), $\|\cdot\|$ denotes the distance between two data sample.

Among the state-of-the-art attacks, C&W [6] attacks achieve the best attacking performance comparing with existing-attack algorithms. This set of attacks provide 3 metrics to measure its distortion ($L_0$, $L_2$ and $L_\infty$ norms),

and $L_2$ norm is the most powerful attack among them. C&W attacks with $L_2$ norm could be formulated as follows,

$$
\begin{aligned}
\min \quad & \|\delta\|_2 + c \cdot Loss(x') \\
s.t. \quad & x' = x + \delta \in D
\end{aligned}
\quad (3)
$$

where $c$ is a hyper-parameter, and *Loss* is defined as

$$Loss(x') = \max(\max\{Z(x')_i : i \neq t\} - Z(x')_t, -k) \quad (4)$$

where $Z(x)$ is the output of the network before the softmax layer, and a hyper-parameter $k$ encourages the attack to search for an adversarial example $x'$ that will be classified as label $t$ with high confidence.

### 3.3. Discrete cosine transform

Discrete cosine transform (DCT) [11,12,46] is a Fourier-like transform, which performs the cosine-series expansion. Due to its high "energy compaction" property, the transformed signal can be easily analyzed using few low-frequency components. Thus, DCT is usually employed to perform decorrelation of the input signal and to present the output in the frequency domain. Among several types of DCT, two-dimensional DCT is the most popular symmetric variation of the transform that operates on images and its inverse. Two-dimensional DCT of an $M \times N$ image $s(x, y)$ are defined as follows.

$$
\begin{aligned}
S(\mu, \nu) = \alpha_\mu \alpha_\nu \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} s(x, y) \\
\times \cos(\frac{\pi\mu(2x+1)}{2M}) \cos(\frac{\pi\nu(2y+1)}{2N})
\end{aligned}
\quad (5)
$$

The inverse of two-dimensional DCT for an $M \times N$ sample:

$$
\begin{aligned}
s(x, y) = \sum_{\mu=0}^{M-1} \sum_{\nu=0}^{N-1} \alpha_\mu \alpha_\nu S(\mu, \nu) \\
\times \cos(\frac{\pi\mu(2x+1)}{2M}) \cos(\frac{\pi\nu(2y+1)}{2N})
\end{aligned}
\quad (6)
$$

The matrix $S(\mu, \nu)$ is the DCT coefficients of image $s(x, y)$, whereas, the basis functions are,

$$
\begin{aligned}
\alpha_\mu = \begin{cases} 1/\sqrt{M} & x = 0 \\ \sqrt{2/M} & 0 \leq x \leq M-1 \end{cases} \\
\alpha_\nu = \begin{cases} 1/\sqrt{N} & y = 0 \\ \sqrt{2/N} & 0 \leq y \leq N-1 \end{cases}
\end{aligned}
\quad (7)
$$

The energy of the image concentrates on the low-frequency DCT components. Human eyes have a poor perception in high-frequency DCT domain. Based on this property, Two-dimensional DCT is widely used in lossy compression [19] and image steganography [10]. In this paper, we apply two-dimensional DCT to availably filter out recessive features from the image.

# 4. Methodology

In this section, we introduce the feature-filter framework and how to design the DCT-based filter to reliably filter out recessive features.

## 4.1. Feature-filter framework

In order to achieve high imperceptibility, state-of-the-art adversarial attacks attempt to enrich recessive features by adding perturbations into original images, which results in a very significant difference in recessive features between benign images and adversarial images. Therefore, adversarial examples is sensitive to the operation of filtering out recessive features, whereas benign examples are not. That is, if the test image is an adversarial example, the DNN-based classifier would predict it and its filtered version into different categories. This idea drove us to design an adversarial example detector by comparing DNN's prediction labels.

Feature-filter employs well-established metamorphic testing to expose adversarial images that lead to inconsistent predictions of the DNN model. As depicted in Fig. 1, the feature-filter consists of a filter, a predictor, and a discriminator, whose functions are shown below.

(1) Filter: Filter out recessive features of the input image and obtain its filtered version;

(2) Predictor: Enter the test image and its filtered version into DNN model $f$ and obtain the DNN's prediction labels;

(3) Discriminator: If the test image and its filtered version are predicted into the same category, the test image is judged to be benign; otherwise, it is judged to be adversarial.

In feature-filter framework, the performance of the proposed detector depends on DNN's prediction on the filtered image. Therefore, a key problem is how to construct a filter so that it can reliably filter out recessive features from the test image. In the next subsection, the construction of the filter is introduced in detail.

## 4.2. DCT-based filter

This section describes how we design a DCT-based filter to filter out the recessive features of the test image. As shown in Fig. 2, the recessive features are usually reserved for high-frequency areas of the image, while the dominant features are reserved for low-frequency areas. Through DCT transform-domain method, we transform spatial-domain image pixels into frequency-domain coefficients and separate the dominant features from the recessive features. By discarding coefficients of high-frequency areas, the filter could filters out most of recessive features. A DCT-based filter is illustrated in Fig. 3. The details of the filter are listed below.

(1) Executing two-dimensional DCT, the $M \times N$ image $x$ is transformed into an $M \times N$ matrix $xf$ of frequency-domain coefficients;

(2) Generating a dominant feature matrix $lf$ by reserving $(\alpha M \times \alpha N)$ low-frequency coefficients and setting all of high-frequency coefficients to 0 in matrix $xf$, where feature reservation ratio $\alpha$ ranges from 0 to 1;

(3) Executing the inverse of two-dimensional DCT, the matrix $lf$ is transformed into a filtered image $xm$, which reserves most of the dominant features and removes a great number of the recessive features.

# 5. Experimental evaluation

## 5.1. Performance of DCT-based filter

The design of the feature-filter roots from an explanation for the existence of adversarial examples with imperceptible perturbations, that is, recessive features are involved in neural network decisions. It is critical to reliably filter out recessive features of the test image. In our work, the dominant features and the recessive features are defined from a human-centric perspective. Therefore, the feature-filter should generate filtered versions which are still clearly recognizable to humans.

To test the DCT-based filter, we set feature reservation ratio $\alpha = 0.5$, which means 3/4 of the frequency features are filtered out by the DCT-based filter. As shown in Fig. 4, five random test images are entered into the DCT-based filter and we obtain their filtered versions. Even if most of the high-frequency features are moved by the DCT-based filter, human eyes still clearly identify test images and filtered versions as the same category. It indicates that the moved features are not perceived by the human eyes.

Whether the DCT-based filter can reliably filter out recessive features from the test image? To verify the performance of the DCT-based filter, we conducted an extensive user study based on human eye evaluation. In the test, the DCT-based filter generates 300 filtered images at multiple feature reservation ratios, which are from 30 benign images in 10 categories in ImageNet dataset [34]. We recruit a total of 50 participants to label these filtered versions. For each trail, filtered images are shown on the screen at a fixed size and the order of images is random. In ImageNet dataset, each image contains more than one object, and only one object category is labeled in each image. In order to clarify ambiguity in the evaluation, we design the following test.

**Test:** Each participant is required to label 300 filtered images generated by the DCT-based filter. The participants are allow to identify multiple (up to 5) objects in an image. As long as one of the objects indeed corresponded to the ground-truth label, this filtered image would be judged to reserve the enough dominant features. We calculate the proportion of trails where filtered images are labeled as a ground-truth label. Higher proportion denotes better per-
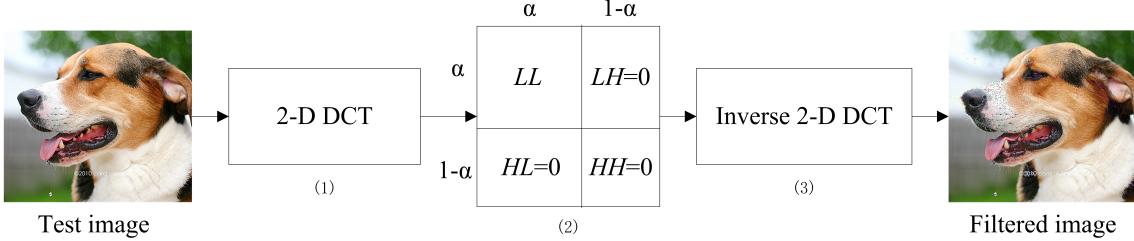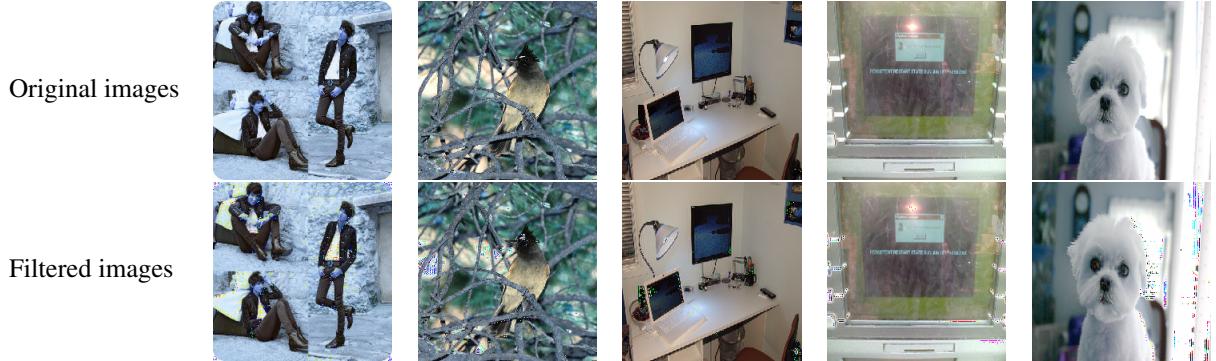
Figure 3. Overview of the DCT-based filter.



Figure 4. Performance of the DCT-based filter.

| $\alpha$ | 0.9 | 0.8 | 0.7 | 0.6 |
|---|---|---|---|---|
| Proportion | **98.73%** | 98.33% | 97.67% | 96.53% |

Table 1. Proportion of filtered images that are correctly labeled by human eyes.

formance of the DCT-based filter.

**Results:** Table 1 lists the proportion of filtered images being labeled correctly at 4 feature reservation ratios $\alpha$. Most of the filtered images are still recognizable to human eyes, which indicates human eyes are not sensitive to the high-frequency features of images. In particular, when $\alpha$ = 0.6, i.e. 64% of the high-frequency features are filtered out, the filtered images are still correctly identified 96.53% of the time. That is to say, human eyes rarely use high-frequency features to understand images. We also observe that the proportion slowly increases as $\alpha$ increases, which denotes the high-frequency features only contain a small number of comprehensible features to humans (dominant features). The results show that the DCT-based filter can reliably filter out recessive features from the test image.

### 5.2. Detection against C&W attack

The design of the feature-filter is based on this assumption that the attacker wants adversarial examples with imperceptible perturbations. To test its detection performance, we need to generate a set of adversarial examples with very imperceptible perturbations. Among state-of-the-art adversarial attacks, C&W attack is proven to be a more powerful attack, in which $L_2$ attack can generate higher quality adversarial examples. C&W attack can bypass multiple types of detectors. Similar results are reported in the literature [5] where C&W attack bypassed all of the ten detection methods. Therefore, we focus on our experiments with the challenging C&W attack ($L_2$ norm) on the CIFAR-10 over Carlini network [6] and ImageNet over Inception V3.

We evaluate the feature-filter from true-positive rate (TPR) and true-negative rate (TNR). TPR measures the proportion of adversarial examples that are correctly identified by the detector, while TNR measures the proportion of benign examples that are correctly identified by the detector. The experiment randomly select 500 CIFAR-10 adversarial examples and 50 ImageNet adversarial examples generated by C&W attack to compute TPR and TNR. Notes that the results exclude the failed adversarial examples and overfitting examples from consideration.

Table 2 lists the results of TPR and TNR at multiple feature reservation ratios $\alpha$. In the best case, the detection rate of adversarial examples is as high as 98.20% for CIFAR-10 when $\alpha$ = 0.8 and 100.00% for ImageNet when $\alpha$ = 0.7. In addition, we also observe that TNR decreases as $\alpha$ decreases, which is consistent with the facts that more feature loss reduces the accuracy of neural networks.

We compare the feature-filter with previous adversarial detection approaches based on image transformation, e.g.,

| $\alpha$ | CIFAR-10 | | ImageNet | |
|---|---|---|---|---|
| | TPR | TNR | TPR | TNR |
| 0.95 | 97.00% | **98.80%** | 72.00% | **100.00%** |
| 0.90 | **98.20%** | 97.00% | 84.00% | 100.00% |
| 0.85 | 97.80% | 96.80% | 96.00% | 100.00% |
| 0.80 | 98.20% | 94.20% | 98.00% | 100.00% |
| 0.75 | 98.00% | 92.60% | 98.00% | 100.00% |
| 0.70 | 98.20% | 89.20% | **100.00%** | 98.00% |
| 0.65 | 97.00% | 83.00% | 100.00% | 98.00% |
| 0.60 | 95.60% | 79.20% | 100.00% | 90.00% |
| 0.55 | 96.60% | 69.60% | 100.00% | 92.00% |
| 0.50 | 96.20% | 66.00% | 100.00% | 90.00% |

Table 2. TPR and TNR at multiple feature reservation ratios $\alpha$.

bit depth reduction [44], spatial smoothing (local smoothing and non-local smoothing) [44], and rotation [42]. Table 3 lists the results of TPR and TNR for several detectors built upon single transformation on the CIFAR-10 dataset. The results show that DCT transform is a more effective approach than other image transformation to distinguish adversarial examples from benign examples.

| Approaches | Parameters | TPR | TNR |
|---|---|---|---|
| Feature-Filter | 0.90 | **98.20%** | **97.00%** |
| | 0.80 | 98.20% | 94.20% |
| | 0.70 | 98.20% | 89.20% |
| | 0.60 | 95.60% | 79.20% |
| | 0.50 | 96.20% | 66.00% |
| Bit Depth Reduction | 1-bit | 89.15% | 45.11% |
| | 2-bit | 92.22% | 78.31% |
| | 3-bit | **93.88%** | 92.69% |
| | 4-bit | 89.00% | 97.52% |
| | 5-bit | 85.75% | **98.76%** |
| Median Smoothing | 2×2 | **95.09%** | **82.03%** |
| | 3×3 | 94.07% | 66.17% |
| | 4×4 | 89.80% | 43.87% |
| Non-local Mean | 11-3-2 | 86.30% | 99.25% |
| | 11-3-4 | **92.73%** | 96.41% |
| | 13-3-2 | 89.72% | **99.26%** |
| | 13-3-4 | 90.63% | 97.03% |
| Rotation | -20 | 88.48% | 60.22% |
| | -10 | 91.60% | 84.01% |
| | 10 | **93.33%** | **85.50%** |
| | 20 | 91.17% | 65.30% |

Table 3. Comparison of TPR and TNR based on several image transformation on CIFAR-10.

We employ the Receiver Operating Characteristic (ROC)

curve and the corresponding Area Under Curve (AUC) to evaluate the detection performance on the CIFAR-10 dataset. Figure 5 plots the ROC curves for several detectors on the task of detecting adversarial images. Table 4 lists the AUC values for several detectors based on image transformation.

The high detection rate shows that filtering out recessive features could significantly change DNN's prediction labels of imperceptible adversarial examples. The results suggest that recessive features of misleading DNNs result in imperceptible adversarial examples, and the adversarial attack enriches these recessive features.
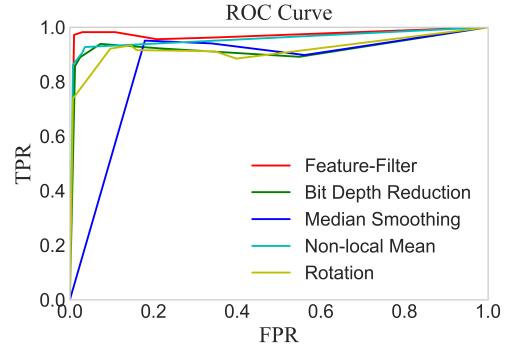


Figure 5. Comparison of ROC curves.

| Approaches | AUC |
|---|---|
| Feature-Filter | **96.23%** |
| Bit Depth Reduction | 92.18% |
| Median Smoothing | 85.70% |
| Non-local Mean | 95.76% |
| Rotation | 91.95% |

Table 4. Comparison of AUC values.

## 5.3. Natural noise detection

Adversarial examples are a kind of images with special noises, while there are a large number of images with natural noise in reality. A good adversarial detector should be able to tolerate natural noises, which would identify natural noise examples to be benign rather than adversarial.

The section verifies that our detector can correctly identify natural noise images as benign ones. We add 3 types of noises, e.g., Gaussian, Poisson, and salt&pepper to generate natural noise images. The natural noise detection performance is evaluated on a pre-trained Inception V3 [40] model in the ImageNet classification task. Notes that natural noises here refer to the perturbations which do not change the prediction of the neural network, so that distinguishing them from adversarial perturbations.

| $\alpha$ | Gaussian | | Poisson | | Salt&Pepper | |
|---|---|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 |
| 0.95 | **79.42%** | **93.31%** | **80.87%** | **94.68%** | **79.29%** | **93.62%** |
| 0.90 | 79.19% | 93.20% | 79.76% | 94.59% | 75.61% | 93.01% |
| 0.85 | 78.08% | 92.93% | 79.50% | 94.18% | 73.98% | 90.94% |
| 0.80 | 79.38% | 92.20% | 79.42% | 93.37% | 74.64% | 89.77% |
| 0.75 | 79.13% | 92.47% | 79.29% | 94.05% | 72.32% | 88.67% |
| 0.70 | 77.72% | 91.53% | 78.38% | 93.21% | 71.78% | 87.92% |
| 0.65 | 77.51% | 91.39% | 77.52% | 92.45% | 71.29% | 88.24% |
| 0.60 | 74.60% | 91.23% | 75.90% | 91.05% | 71.33% | 88.79% |
| 0.55 | 73.08% | 89.76% | 75.67% | 90.85% | 70.05% | 88.28% |
| 0.50 | 70.79% | 88.62% | 72.53% | 89.66% | 67.62% | 86.19% |

Table 5. Natural noise detection on the ImageNet.

Table 5 lists the accuracy of the feature-filter in detecting 3 types of natural noise images. Due to multiple objects in an image of the ImageNet dataset, We test top-1 and top-5 to evaluate our detector performance on natural noise images. Higher accuracy denotes better tolerance of the detector to the natural noise. As shown in Tab. 5, the top-5 accuracy is significantly higher than that of the top-1. Particularly, the top-5 accuracy of the feature-filter is close to 95% for Poisson noise. And our detector shows an approximate performance for 3 types of natural noises, which can stably avoid the interference of natural noise. We also see that the accuracy of the feature-filter slowly decreases as $\alpha$ decreases. The phenomenon reveals that high-frequency features only contain a small amount of dominant features. Therefore, the DCT-based filter could be regarded as a reliable tool to filter out the recessive features from the test image.

### 5.4. Computation cost

Apart from detection rate, efficiency is also an important parameter for a good adversarial detector. The feature-filter consists of a filter, a predictor and a discriminator. The filter is based on two-dimensional DCT to generate the filtered image of the test image. As defined in Eqs. 5 and 6, the time complexity of DCT and its inverse transform is $O(M \times N)$ for an $M \times N$ image. The predictor is made up of the target DNN, which performs two predictions on the test image and its filtered version. Thus, computation cost of the predictor depends on the efficiency of the neural network. The discriminator only gives the decision according to two predicted labels by the predictor, thus its time complexity is $O(1)$.

The actual running time depends on many factors including experimental setup, programming skill, etc. Our test is implemented 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). We choose 1,000 random images to test the running time of the feature-filter. The average running time is 0.0351 second for a $224 \times 224 \times 3$ ImageNet image and 0.0028 second for a $32 \times 32 \times 3$ CIFAR-10 image. Due to the low computation cost, the feature-filter is considered a real-time adversarial detector.

## 6. Conclusion

The explanation of adversarial examples remains an open issue. As a human-centric phenomenon, we divide the image feature into dominant features and recessive features. We reveal that the existence of imperceptible adversarial examples is due to the DNN's use of the recessive features that mislead neural networks. Adversarial attacks attempt to enrich these recessive features so that the perturbed example can "fool" the DNN without noticeable artifacts.

According to above explanation, adversarial examples should be more sensitive to the operation of filtering out recessive features than benign/noisy examples. Inspired by this idea, we propose an adversarial detection approach named feature-filter. The feature-filter determines whether the test image is adversarial or benign by comparing the DNN's prediction labels on the test input and its filtered version. To filter out recessive features, we introduce two-dimensional DCT to separate the dominant features and the recessive features. The experimental results show that the DCT-based filter can reliably filter out the recessive features of the image, and the feature-filter is a label-only tool to detect adversarial examples in real time.

The effectiveness of the feature-filter stems from a reliable filtering of recessive features. As discussed in Sec. 5, two-dimensional DCT is an approximate method to separate the recessive features from the dominant features. Future works are to discover more reliable feature filtering methods or combine the feature-filter with other methods, so that the adversarial detector is more accurate and robust.

# References

[1] Yuval Bahat, Michal Irani, and Gregory Shakhnarovich. Natural and adversarial error detection using invariance to image transformations. *CoRR*, abs/1902.00236, 2019. 3

[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, 4

[3] Ramon Blanco-Gonzalo, Oscar Miguel-Hurtado, Chiara Lunerti, Richard M. Guest, Barbara Corsetti, Elakkiya Ellavarason, and Raul Sánchez-Reillo. Biometric systems interaction assessment: The state of the art. *IEEE Trans. Hum. Mach. Syst.*, 49(5):397–410, 2019. 1

[4] Sébastien Bubeck, Eric Price, and Ilya P. Razenshteyn. Adversarial examples from computational constraints. *CoRR*, abs/1805.10204, 2018. 2

[5] Nicholas Carlini and David A. Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In Bhavani M. Thuraisingham, Battista Biggio, David Mandell Freeman, Brad Miller, and Arunesh Sinha, editors, *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec@CCS 2017, Dallas, TX, USA, November 3, 2017*, pages 3–14. ACM, 2017. 6

[6] 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, 3, 4, 6

[7] Florent Chiaroni, Mohamed-Cherif Rahal, Nicolas Hueber, and Frédéric Dufaux. Self-supervised learning for autonomous vehicles perception: A conciliation between analytical and learning methods. *IEEE Signal Process. Mag.*, 38(1):31–41, 2021. 1

[8] Gilad Cohen, Guillermo Sapiro, and Raja Giryes. Detecting adversarial samples using influence functions and nearest neighbors. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 14441–14450. IEEE, 2020. 1

[9] Gilad Cohen, Guillermo Sapiro, and Raja Giryes. Detecting adversarial samples using influence functions and nearest neighbors. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 14441–14450. IEEE, 2020. 3

[10] Hong-Zhu Dai, Jie Cheng, and Yafeng Li. A novel steganography algorithm based on quantization table modification and image scrambling in DCT domain. *Int. J. Pattern Recognit. Artif. Intell.*, 35(1):2154001:1–2154001:18, 2021. 4

[11] Nilaksh Das, Madhuri Shanbhogue, Shang-Tse Chen, Fred Hohman, Li Chen, Michael E. Kounavis, and Duen Horng Chau. Keeping the bad guys out: Protecting and vaccinating deep learning with JPEG compression. *CoRR*, abs/1705.02900, 2017. 2, 4

[12] Gintare Karolina Dziugaite, Zoubin Ghahramani, and Daniel M. Roy. A study of the effect of JPG compression on adversarial images. *CoRR*, abs/1608.00853, 2016. 2, 4

[13] Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner. Detecting adversarial samples from artifacts. *CoRR*, abs/1703.00410, 2017. 3

[14] 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, 3, 4

[15] 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, 4

[16] 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. IEEE Computer Society, 2018. 1, 4

[17] 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, 4

[18] 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. 2

[19] Yiannis Kantaros, Taylor J. Carpenter, Sangdon Park, Radoslav Ivanov, Sooyong Jang, Insup Lee, and James Weimer. Visionguard: Runtime detection of adversarial inputs to perception systems. *CoRR*, abs/2002.09792, 2020. 3, 4

[20] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. 1

[21] Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense against unseen threat models. *CoRR*, abs/2006.12655, 2020. 1, 4

[22] 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, 4

[23] Giulio Lovisotto, Henry Turner, Ivo Sluganovic, Martin Strohmeier, and Ivan Martinovic. SLAP: improving physical adversarial examples with short-lived adversarial perturbations. In *30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021*, pages 1865–1882, 2021. 1

[24] Shiqing Ma, Yingqi Liu, Guanhong Tao, Wen-Chuan Lee, and Xiangyu Zhang. NIC: detecting adversarial samples with neural network invariant checking. In *26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019*. The Internet Society, 2019. 3

[25] Saeed Mahloujifar, Dimitrios I. Diochnos, and Mohammad Mahmoody. The curse of concentration in robust learning: Evasion and poisoning attacks from concentration of measure. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 4536–4543. AAAI Press, 2019. 2

[26] Yisroel Mirsky, Tomer Doitshman, Yuval Elovici, and Asaf Shabtai. Kitsune: An ensemble of autoencoders for online network intrusion detection. *CoRR*, abs/1802.09089, 2018. 1

[27] 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, 3, 4

[28] 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, 3, 4

[29] Nicolas Papernot, Patrick D. McDaniel, Arunesh Sinha, and Michael P. Wellman. Towards the science of security and privacy in machine learning. *CoRR*, abs/1611.03814, 2016. 3

[30] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy, SP 2016, San Jose, CA, USA, May 22-26, 2016*, pages 582–597. IEEE Computer Society, 2016. 1

[31] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy, SP 2016, San Jose, CA, USA, May 22-26, 2016*, pages 582–597. IEEE Computer Society, 2016. 3

[32] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. Deepxplore: Automated whitebox testing of deep learning systems. In *Proceedings of the 26th Symposium on Operating Systems Principles, Shanghai, China, October 28-31, 2017*, pages 1–18. ACM, 2017. 1

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

*USA, June 26 - July 1, 2016*, pages 410–417. IEEE Computer Society, 2016. 3

[34] 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 Fei-Fei Li. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.*, 115(3):211–252, 2015. 5

[35] 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. 2

[36] Uri Shaham, Yutaro Yamada, and Sahand Negahban. Understanding adversarial training: Increasing local stability of supervised models through robust optimization. *Neurocomputing*, 307:195–204, 2018. 1

[37] 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, 4

[38] Philip Sperl, Ching-Yu Kao, Peng Chen, Xiao Lei, and Konstantin Böttinger. DLA: dense-layer-analysis for adversarial example detection. In *IEEE European Symposium on Security and Privacy, EuroS&P 2020, Genoa, Italy, September 7-11, 2020*, pages 198–215. IEEE, 2020. 3

[39] 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, 4

[40] 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. 4, 7

[41] 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, 3

[42] Shixin Tian, Guolei Yang, and Ying Cai. Detecting adversarial examples through image transformation. 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 4139–4146. AAAI Press, 2018. 3, 4, 7

[43] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 1

[44] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. In *25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018*. The Internet Society, 2018. 3, 7

[45] Xing Xu, Jingran Zhang, Yujie Li, Yichuan Wang, Yang Yang, and Heng Tao Shen. Adversarial attack against urban scene segmentation for autonomous vehicles. *IEEE Trans. Ind. Informatics*, 17(6):4117–4126, 2021. 1

[46] Zakia Yahya, Muhammad Hassan, Muhammad Shahzad Younis, and Muhammad Shafique. Probabilistic analysis of targeted attacks using transform-domain adversarial examples. *IEEE Access*, 8:33855–33869, 2020. 2, 4

[47] Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael I. Jordan. ML-LOO: detecting adversarial examples with feature attribution. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 6639–6647. AAAI Press, 2020. 3

[48] Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. Adversarial examples: Attacks and defenses for deep learning. *IEEE Trans. Neural Networks Learn. Syst.*, 30(9):2805–2824, 2019. 1, 4