

Feature Squeezing: Detecting Adversarial Examples in Deep Neural Networks

Weilin Xu, David Evans, Yanjun Qi
University of Virginia
evadeML.org

Abstract—Although deep neural networks (DNNs) have achieved great success in many tasks, recent studies have shown they are vulnerable to adversarial examples. Such examples, typically generated by adding small but purposeful distortions, can frequently fool DNN models. Previous studies to defend against adversarial examples mostly focused on refining the DNN models, but have either shown limited success or suffered from expensive computation. We propose a new strategy, *feature squeezing*, that can be used to harden DNN models by detecting adversarial examples. Feature squeezing reduces the search space available to an adversary by coalescing samples that correspond to many different feature vectors in the original space into a single sample. By comparing a DNN model’s prediction on the original input with that on squeezed inputs, feature squeezing detects adversarial examples with high accuracy and few false positives. This paper explores two types of feature squeezing: reducing the color bit depth of each pixel and spatial smoothing. These strategies are inexpensive and complementary to other defenses, and can be combined in a joint detection framework to achieve high detection rates against state-of-the-art attacks.

I. INTRODUCTION

Deep Neural Networks (DNNs) perform exceptionally well on many artificial intelligence tasks, including security-sensitive applications like malware classification [26], [8] and face recognition [35]. Unlike when machine learning is used in other fields, security applications involve intelligent and adaptive adversaries responding to the deployed systems. Recent studies have shown that attackers can force deep learning object classification models to mis-classify images by making imperceptible modifications to pixel values. The maliciously generated inputs are called “adversarial examples” [10], [39] and are normally crafted using an optimization procedure to search for small, but effective, artificial perturbations.

The goal of this work is to harden DNN systems against adversarial examples by detecting them successfully. Detecting an attempted attack may be as important as predicting correct outputs. When running locally, a classifier that can detect adversarial inputs may alert its users or take fail-safe actions (e.g., a fully autonomous drone returns to its base) when it spots adversarial inputs. For an on-line classifier whose model is being used (and possibly updated) through API calls from external clients, the ability to detect adversarial examples may enable the operator to identify malicious clients and exclude their inputs. Another reason that detecting adversarial examples is important is because even with the strongest defenses, adversaries will occasionally be able to get lucky and find an adversarial input. For asymmetrical security applications like malware detection, the adversary may only need to find a single example that preserves the desired malicious behavior but is classified as benign to launch a successful attack. This seems like a hopeless situation for an on-line classifier operator, but

the game changes if the operator can detect even unsuccessful attempts during an adversary’s search process.

Most of the previous work aiming to harden DNN systems, including like *adversarial training* and *gradient masking* (details in Section II-C), focused on modifying the DNN models themselves. In contrast, our work focuses on finding simple and low-cost defensive strategies that alter the input samples but leave the model unchanged. A few other recent studies have proposed methods to detect adversarial examples through sample statistics, training a detector, or prediction inconsistency (Section II-D). Our approach, which we call *feature squeezing*, is driven by the observation that the feature input spaces are often unnecessarily large, and this vast input space provides extensive opportunities for an adversary to construct adversarial examples. Our strategy is to reduce the degrees of freedom available to an adversary by “squeezing” out unnecessary input features.

The key to our approach is to compare the model’s prediction on the original sample with its prediction on the sample after squeezing, as depicted in Figure 1. If the original and squeezed inputs produce substantially different outputs from the model, the input is likely to be adversarial. By comparing the difference between predictions with a selected threshold value, our system outputs the correct prediction for legitimate examples and rejects adversarial inputs.

The approach generalizes to other domains where deep learning is used, such as voice recognition and natural language processing. Carlini et al. have demonstrated that lowering the sampling rate helps to defend against the adversarial voice commands [4]. Hosseini et al. proposed to perform spell checking on the inputs of a character-based toxic text detection system to defend against the adversarial examples [16]. Both of them could be regard as an instance of feature squeezing.

Although feature squeezing generalizes to other domains, here we focus on image classification because it is the domain where adversarial examples have been most extensively stud-

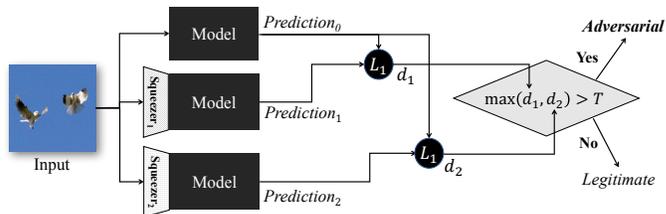


Fig. 1: Detecting adversarial examples. The model is evaluated on both the original input and the input after being pre-processed by one or more feature squeezers. If any of the predictions on the squeezed inputs are too different from the original prediction, the input is determined to be adversarial.

ied. We explore two simple methods for squeezing features of images: reducing the color depth of each pixel in an image, and using spatial smoothing to reduce the differences among individual pixels. We demonstrate that feature squeezing significantly enhances the robustness of a model by predicting correct labels of adversarial examples, while preserving the accuracy on legitimate inputs (Section IV), thus enabling an accurate detector for adversarial examples (Section V). Feature squeezing appears to be both more accurate and general, and less expensive, than previous methods.

Contributions. Our key contribution is introducing and evaluating feature squeezing as a technique for detecting adversarial examples. We introduce the general detection framework (depicted in Figure 1), and show how it can be instantiated to accurately detect adversarial examples generated by a wide range of state-of-the-art methods.

We study two instances of feature squeezing: reducing color bit depth (Section III-A) and both local and non-local spatial smoothing (Section III-B). We report on experiments that show feature squeezing helps DNN models predict correct classification on adversarial examples generated by eleven different and state-of-the-art attacks (Section IV).

Section V explains how we use feature squeezing for detecting adversarial inputs in two distinct situations. In the first case, we (overly-optimistically) assume the model operator knows the attack type and can select a single squeezer for detection. Our results show that the effectiveness of different squeezers against various attacks varies. For instance, the 1-bit depth reduction squeezer achieves a perfect 100% detection rate on MNIST for six different attacks. However, this squeezer is not as effective against those attacks making substantial changes to a small number of pixels (that can be detected well by median smoothing). The model operator normally does not know what attacks an adversary may use, so requires a detection system to work well against any attack. We propose combining multiple squeezers in a joint detection framework. Our experiments show that joint-detection can successfully detect adversarial examples from eleven state-of-the-art attacks at the detection rates of 98% on MNIST and 85% on CIFAR-10 and ImageNet, with low (below 5%) false positive rates.

Feature squeezing is complementary to other adversarial defenses since it does not change the underlying model, and can readily be composed with other defenses such as adversarial training (Section IV-E). Although we cannot guarantee an adaptive attacker cannot succeed against a particular feature squeezing configuration, our results show it is effective against state-of-the-art methods, and it considerably complicates the task of an adaptive adversary even with full knowledge of the model and defense (Section V-D).

II. BACKGROUND

This section provides a brief introduction to neural networks, methods for finding adversarial examples, and previously-proposed defenses.

A. Neural Networks

Deep Neural Networks (DNNs) can efficiently learn highly-accurate models from large corpora of training samples in

many domains [19], [13], [26]. Convolutional Neural Networks (CNNs), first popularized by LeCun et al. [21], perform exceptionally well on image classification. A deep CNN can be written as a function $g : X \rightarrow Y$, where X represents the input space and Y is the output space representing a categorical set. For a sample, $\mathbf{x} \in X$,

$$g(\mathbf{x}) = f_L(f_{L-1}(\dots((f_1(\mathbf{x}))))).$$

Each f_i represents a layer, which can be a classical feed-forward linear layer, rectification layer, max-pooling layer, or a convolutional layer that performs a sliding window operation across all positions in an input sample. The last output layer, f_L , learns the mapping from a hidden space to the output space (class labels) through a softmax function.

A training set contains N_{tr} labeled inputs in which the i -th input is denoted (\mathbf{x}_i, y_i) . When training a deep model, parameters related to each layer are randomly initialized, and input samples (\mathbf{x}_i, y_i) are fed through the network. The output of this network is a prediction $g(\mathbf{x}_i)$ associated with the i -th sample. To train the DNN, the difference between prediction output, $g(\mathbf{x}_i)$, and its true label, y_i , is fed back into the network using a back-propagation algorithm to update DNN parameters.

B. Generating Adversarial Examples

An adversarial example is an input crafted by an adversary with the goal of producing an incorrect output from a target classifier. Since ground truth, at least for image classification tasks, is based on human perception which is hard to model or test, research in adversarial examples typically defines an adversarial example as a misclassified sample \mathbf{x}' generated by perturbing a correctly-classified sample \mathbf{x} (*a.k.a* seed example) by some limited amount.

Adversarial examples can be *targeted*, in which case the adversary’s goal is for \mathbf{x}' to be classified as a particular class t , or *untargeted*, in which case the adversary’s goal is just for \mathbf{x}' to be classified as any class other than its correct class. More formally, given $\mathbf{x} \in X$ and $g(\cdot)$, the goal of an targeted adversary with target $t \in Y$ is to find an $\mathbf{x}' \in X$ such that

$$g(\mathbf{x}') = t \wedge \Delta(\mathbf{x}, \mathbf{x}') \leq \epsilon \quad (1)$$

where $\Delta(\mathbf{x}, \mathbf{x}')$ represents the difference between input \mathbf{x} and \mathbf{x}' . An untargeted adversary’s goal is to find an $\mathbf{x}' \in X$ such that

$$g(\mathbf{x}') \neq g(\mathbf{x}) \wedge \Delta(\mathbf{x}, \mathbf{x}') \leq \epsilon. \quad (2)$$

The strength of the adversary, ϵ , measures the permissible transformations. The distance metric, $\Delta(\cdot)$, and the adversarial strength threshold, ϵ , are meant to model how close an adversarial example \mathbf{x}' needs to be to the original sample \mathbf{x} to “fool” a human observer.

Several techniques have been proposed to find adversarial examples. Szegedy et al. [39] first observed that DNN models are vulnerable to adversarial perturbation and used the Limited-memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) algorithm to find adversarial examples. Their study also found that adversarial perturbations generated from one DNN model can also force other DNN models to produce incorrect outputs. Subsequent papers have explored other strategies to

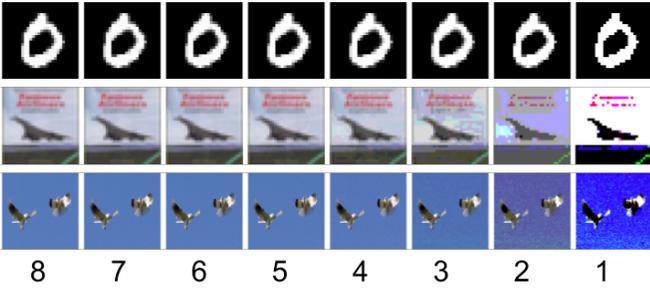


Fig. 2: Image examples with bit depth reduction. The first column shows images from MNIST, CIFAR-10 and ImageNet, respectively. Other columns show squeezed versions at different color-bit depths, ranging from 8 (original) to 1.

generate adversarial manipulations, including using the linear assumption behind a model [10], [28], saliency maps [32], and evolutionary algorithms [29].

Equations (1) and (2) suggest two different parameters for categorizing methods for finding adversarial examples: whether they are targeted or untargeted, and the choice of $\Delta()$, which is typically an L_p -norm distance metric. When given a m -dimensional vector $\mathbf{z} = \mathbf{x} - \mathbf{x}' = (z_1, z_2, \dots, z_m)^T \in \mathbb{R}^m$, the L_p norm is defined by:

$$\|\mathbf{z}\|_p = \sqrt[p]{\sum_{i=1}^m |z_i|^p} \quad (3)$$

The three norms used as $\Delta()$ choices for popular adversarial methods are:

- L_∞ : $\|\mathbf{z}\|_\infty = \max_i |z_i|$. The L_∞ norm measures the maximum change in any dimension. This means an L_∞ attack is limited by the maximum change it can make to each pixel, but can alter all the pixels in the image by up to that amount.
- L_2 : $\|\mathbf{z}\|_2 = \sqrt{\sum_i z_i^2}$. The L_2 norm corresponds to the Euclidean distance between \mathbf{x} and \mathbf{x}' . This distance can remain small when many small changes are applied to many pixels.
- L_0 : $\|\mathbf{z}\|_0 = \#\{i \mid z_i \neq 0\}$. For images, this metric measures the number of pixels that have been altered between \mathbf{x} and \mathbf{x}' , so an L_0 attack is limited by the number of pixels it can alter.

We discuss the eleven attacking algorithms, grouped by the norm they used for Δ , used in our experiments further below.

1) Fast Gradient Sign Method: FGSM (L_∞ , Untargeted)

Goodfellow et al. hypothesized that DNNs are vulnerable to adversarial perturbations because of their linear nature [10]. They proposed the *fast gradient sign method* (FGSM) for efficiently finding adversarial examples. To control the cost of attacking, FGSM assumes that the attack strength at every feature dimension is the same, essentially measuring the perturbation $\Delta(\mathbf{x}, \mathbf{x}')$ using the L_∞ -norm. The strength of perturbation at every dimension is limited by the same constant parameter, ϵ , which is also used as the amount of perturbation.

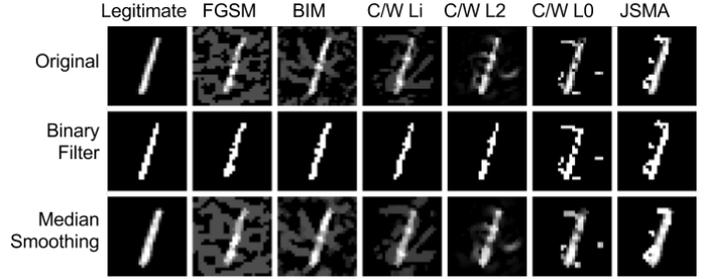


Fig. 3: Examples of adversarial attacks and feature squeezing methods extracted from the MNIST dataset. The first column shows the original image and its squeezed versions, while the other columns present the adversarial variants. All targeted attacks are targeted-next.

As an untargeted attack, the perturbation is calculated directly by using gradient vector of a loss function:

$$\Delta(\mathbf{x}, \mathbf{x}') = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(g(\mathbf{x}), y)) \quad (4)$$

Here the loss function, $J(\cdot, \cdot)$, is the loss that have been used in training the specific DNN model, and y is the correct label for \mathbf{x} . Equation (4) essentially increases the loss $J(\cdot, \cdot)$ by perturbing the input \mathbf{x} based on a transformed gradient.

2) Basic Iterative Method: BIM (L_∞ , Untargeted)

Kurakin et al. extended the FGSM method by applying it multiple times with small step size [20]. This method clips pixel values of intermediate results after each step to ensure that they are in an ϵ -neighborhood of the original image \mathbf{x} . For the m -th iteration,

$$\mathbf{x}'_{m+1} = \mathbf{x}'_m + \text{Clip}_{\mathbf{x}, \epsilon} \{\alpha \cdot \text{sign}(\nabla_{\mathbf{x}} J(g(\mathbf{x}'_m), y))\} \quad (5)$$

The clipping equation, $\text{Clip}_{\mathbf{x}, \epsilon}(\mathbf{z})$, performs per-pixel clipping on \mathbf{z} so the result will be in the L_∞ ϵ -neighborhood of the source \mathbf{x} [20].

3) DeepFool (L_2 , Untargeted)

Moosavi et al. used a L_2 minimization-based formulation, termed DeepFool, to search for adversarial examples [28]:

$$\Delta(\mathbf{x}, \mathbf{x}') := \arg \min_{\mathbf{z}} \|\mathbf{z}\|_2, \text{ subject to: } g(\mathbf{x} + \mathbf{z}) \neq g(\mathbf{x}) \quad (6)$$

DeepFool searches for the minimal perturbation to fool a classifier and uses concepts from geometry to direct the search. For linear classifiers (whose decision boundaries are linear planes), the region of the space describing a classifier's output can be represented by a polyhedron (whose plane faces are those boundary planes defined by the classifier). Then DeepFool searches within this polyhedron for the minimal perturbation that can change the classifiers decision. For general non-linear classifiers, this algorithm uses an iterative linearization procedure to get an approximated polyhedron.

4) Jacobian Saliency Map Approach: JSMA (L_0 , Targeted)

Papernot et al. [32] proposed the *Jacobian-based saliency map approach* (JSMA) to search for adversarial examples by only modifying a limited number of input pixels in an image. As a targeted attack, JSMA iteratively perturbs pixels

in an input image that have high adversarial saliency scores. The adversarial saliency map is calculated from the Jacobian (gradient) matrix $\nabla_{\mathbf{x}}g(\mathbf{x})$ of the DNN model $g(\mathbf{x})$ at the current input \mathbf{x} . The $(c, p)^{\text{th}}$ component in Jacobian matrix $\nabla_{\mathbf{x}}g(\mathbf{x})$ describes the derivative of output class c with respect to feature pixel p . The adversarial saliency score of each pixel is calculated to reflect how this pixel will increase the output score of the target class t versus changing the score of all other possible output classes. The process is repeated until classification into the target class is achieved, or it reaches the maximum number of perturbed pixels. Essentially, JSMA optimizes Equation (2) by measuring perturbation $\Delta(\mathbf{x}, \mathbf{x}')$ through the L_0 -norm.

5) Carlini/Wagner Attacks (L_2 , L_∞ and L_0 , Targeted)

Carlini and Wagner recently introduced three new gradient-based attack algorithms that are more effective than all previously-known methods in terms of the adversarial success rates achieved with minimal perturbation amounts [6]. There are versions of their attacks for L_2 , L_∞ , and L_0 norms.

The CW_2 attack formalizes the task of generating adversarial examples as an optimization problem with two terms as usual: the prediction term and the distance term. However, it makes the optimization problem easier to solve with several techniques. The first is using the logits-based objective function instead of the softmax-cross-entropy loss that is commonly used in other optimization-based attacks. This makes it robust against the defensive distillation method [34]. The second is converting the target variable to the *argtanh* space to bypass the box-constraint on the input, making it more flexible in taking advantage of modern optimization solvers, such as Adam. It also uses a binary search algorithm to select a suitable coefficient that performs a good trade-off between the prediction and the distance terms. These improvements enable the CW_2 attack to find adversarial examples with smaller perturbations than previous attacks.

Their CW_∞ attack recognizes the fact that L_∞ norm is hard to optimize and only the maximum term is penalized. Thus, it revises the objective into limiting perturbations to be less than a threshold τ (initially 1, decreasing in each iteration). The optimization reduces τ iteratively until no solution can be found. Consequently, the resulting solution has all the perturbations smaller than the specified τ .

The basic idea of the CW_0 attack is to iteratively use CW_2 to find the least important features and freeze them (so value will never be changed) until the L_2 attack fails with too many features being frozen. As a result, only those features with significant impact on the prediction are changed. This is the opposite of JSMA, which iteratively selects the most important features and performs large perturbations until it successfully fools the target classifier.

C. Defensive Techniques

Papernot et al. [33] provide a comprehensive summary of work on defending against adversarial samples, grouping work into two broad categories: *adversarial training* and *gradient masking*, which we discuss further below. A third approach is to modify feature sets, but it has not previously been applied to DNN models. Wang et al. proposed a theory that unnecessary

features are the primary cause of a classifier’s vulnerability to adversarial examples [41]. Zhang et al. proposed an adversary-aware feature selection model that can improve classifier robustness against evasion attacks [43]. Our proposed feature squeezing method is broadly part of this theme.

Adversarial Training. *Adversarial training* introduces discovered adversarial examples and the corresponding ground truth labels to the training dataset [10], [39]. Ideally, the model will learn how to restore the ground truth from the adversarial perturbations and perform robustly on the future adversarial examples. This technique, however, suffers from the high cost to generate adversarial examples and (at least) doubles the training cost of DNN models due to its iterative re-training procedure. Its effectiveness also depends on having a technique for efficiently generating adversarial examples similar to the one used by the adversary, which may not be the case in practice. As pointed out by Papernot et al. [33], it is essential to include adversarial examples produced by all known attacks in adversarial training, since this defensive training is non-adaptive. But, it is computationally expensive to find adversarial inputs by most known techniques, and there is no way to be confident the adversary is limited to techniques that are known to the trainer.

Gradient Masking. These defenses seek to reduce the sensitivity of DNN models to small changes made to their sample inputs, by forcing the model to produce near-zero gradients. Gu et al. proposed adding a gradient penalty term in the objective function, which is defined as the summation of the layer-by-layer Frobenius norm of the Jacobian matrix [12]. Although the trained model behaves more robustly against adversaries, the penalty significantly reduces the capacity of the model and sacrifices accuracy on many tasks [33]. Papernot et al. introduced defensive distillation to harden DNN models [34]. A defensively distilled model is trained with the smoothed labels generated by a normally-trained DNN model. Then, to hide model’s gradient information from an adversary, the distilled model replaces its last layer with a “harder” softmax function after training. Experimental results found that larger perturbations are required when using JSMA to evade distilled models. However, two subsequent studies showed that defensive distillation failed to mitigate a variant of JSMA with a division trick [5] and a black-box attack [31]. Papernot et al. concluded that methods designed to conceal gradient information are bound to have limited success because of the transferability of adversarial examples [33].

D. Detecting Adversarial Examples

A few recent studies [25], [11], [9] have focused on detecting adversarial examples. The strategies they explored can be considered into three groups: *sample statistics*, *training a detector* and *prediction inconsistency*.

Sample Statistics. Grosse et al. [11] propose a statistical test method for detecting adversarial examples using maximum mean discrepancy and energy distance as the statistical distance measures. Their method requires a large set of adversarial examples and legitimate samples and is not capable of detecting individual adversarial examples, making it less useful in practice. Feinman et al. propose detecting adversarial examples us-

ing kernel density estimation [9], which measures the distance between an unknown input example and a group of legitimate examples in a manifold space (represented as features in some middle layers of a DNN). It is computationally expensive and can only detect adversarial examples lying far from the manifolds of the legitimate population. Using sample statistics to differentiate between adversarial examples and legitimate inputs seems unlikely to be effective against broad classes of attacks due to the intrinsically deceptive nature of such examples. Experimental results from both Grosse et al. [11] and Feinman et al. [9] have found that strategies relying on sample statistics gave inferior detection performance compared to other strategies.

Training a Detector. Similar to adversarial training, adversarial examples can also be used to train a detector. Because of the large number of adversarial examples needed, this method is expensive and prone to overfitting employed adversarial techniques. Metzen et al. proposed attaching a CNN-based detector as a branch off a middle layer of the original DNN [25]. The detector outputs two classes and uses adversarial examples (as one class) plus legitimate examples (as the other class) for training. The detector is trained while freezing the weights of the original DNN, so does not sacrifice classification accuracy on the legitimate inputs. Grosse et al. demonstrate a similar detection method (previously proposed by Nguyen et al. [29]) that adds a new “adversarial” class in the last layer of the DNN model [11]. The revised model is trained with both legitimate and adversarial inputs, reducing the accuracy on legitimate inputs due to the change to the model architecture.

Prediction Inconsistency. The basic idea of *prediction inconsistency* is to measure the disagreement among several models in predicting an unknown input example, since one adversarial example may not fool every DNN model. Feinman et al. borrowed an idea from dropout [15] and designed a detection technique they called *Bayesian neural network uncertainty* [9]. In its original form, a dropout layer randomly drops some weights (by temporarily setting to zero) in each training iteration and uses all weights at the testing phase, which can be interpreted as training many different sub-models and averaging their predictions in testing. For detecting adversarial examples, Feinman et al. propose using the “training” mode of dropout layers to generate many predictions of each input. They reported that the disagreement among the predictions of sub-models is rare on legitimate inputs but common on adversarial examples, thus can be employed for detection.

III. FEATURE SQUEEZING METHODS

Although the notion of feature squeezing is quite general, we focus on two simple types of squeezing: reducing the color depth of images (Section III-A), and using smoothing (both local and non-local) to reduce the variation among pixels (Section III-B). Section IV looks at the impact of each squeezing method on classifier accuracy and robustness against adversarial inputs. These results enable feature squeezing to be used for detecting adversarial examples in Section V.

A. Color Depth

A neural network, as a differentiable model, assumes that the input space is continuous. However, digital computers only

support discrete representations as approximations of continuous natural data. A standard digital image is represented by an array of pixels, each of which is usually represented as a number that represents a specific color.

Common image representations use color bit depths that lead to irrelevant features, so we hypothesize that reducing bit depth can reduce adversarial opportunity without harming classifier accuracy. Two common representations, which we focus on here because of their use in our test datasets, are 8-bit grayscale and 24-bit color. A grayscale image provides $2^8 = 256$ possible values for each pixel. An 8-bit value represents the intensity of a pixel where 0 is black, 255 is white, and intermediate numbers represent different shades of gray. The 8-bit scale can be extended to display color images with separate red, green and blue color channels. This provides 24 bits for each pixel, representing $2^{24} \approx 16$ million different colors.

1) Squeezing Color Bits

While people usually prefer larger bit depth as it makes the displayed image closer to the natural image, large color depths are often not necessary for interpreting images (for example, people have no problem recognizing most black-and-white images). We investigate the bit depth squeezing with three popular datasets for image classification: MNIST, CIFAR-10 and ImageNet.

Greyscale Images (MNIST). The MNIST dataset contains 70,000 images of hand-written digits (0 to 9). Of these, 60,000 images are used as training data and the remaining 10,000 images are used for testing. Each image is 28×28 pixels, and each pixel is encoded as 8-bit grayscale.

Figure 2 shows one example of class 0 in the MNIST dataset in the first row, with the original 8-bit grayscale images in the leftmost and the 1-bit monochrome images rightmost. The rightmost images, generated by applying a binary filter with 0.5 as the cutoff, appear nearly identical to the original images on the far left. The processed images are still recognizable to humans, even though the feature space is only $1/128^{\text{th}}$ the size of the original 8-bit grayscale space.

Figure 3 hints at why reducing color depth can mitigate adversarial examples generated by multiple attack techniques. The top row shows one original example of class 1 from the MNIST test set and six different adversarial examples. The middle row shows those examples after reducing the bit depth of each pixel into binary. To a human eye, the binary-filtered images look more like the correct class; in our experiments, we find this is true for DNN classifiers also (Table III in Section IV).

Color Images (CIFAR-10 and ImageNet). We use two datasets of color images in this paper: the CIFAR-10 dataset with tiny images and the ImageNet dataset with high-resolution photographs. The CIFAR-10 dataset contains 60,000 images, each with 32×32 pixels encoded with 24-bit color and belonging to 10 different classes. The ImageNet dataset is provided by ImageNet Large Scale Visual Recognition Challenge 2012 for the classification task, which contains 1.2 million training images and the other 50,000 images for validation. The photographs in the ImageNet dataset are in different sizes

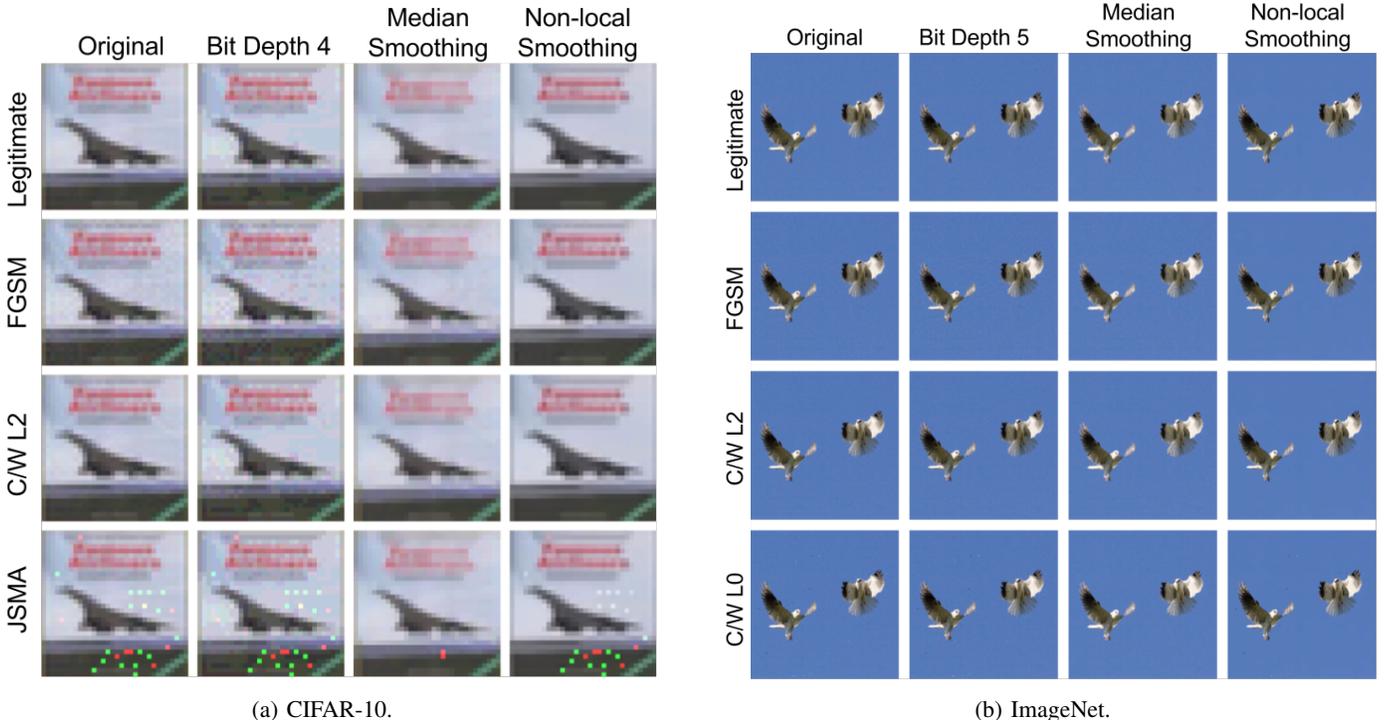


Fig. 4: Examples of adversarial attacks and feature squeezing methods extracted from the CIFAR-10 and ImageNet datasets. The first row presents the original image and its squeezed versions, while the other rows presents the adversarial variants.

and hand-labeled with 1,000 classes. However, they are pre-processed to 224×224 pixels encoded with 24-bit True Color for the target model MobileNet [17], [24] we use in this paper.

The middle row and the bottom row of Figure 2 show that we can reduce the original 8-bit (per RGB channel) images to fewer bits without significantly decreasing the image recognizability to humans. It is difficult to tell the difference between the original images with 8-bit per channel color and images using as few as 4 bits of color depth. Unlike what we observed in the MNIST dataset, however, bit depths lower than 4 do introduce some human-observable loss. This is because we lose much more information in the color image even though we reduce to the same number of bits per channel. For example, if we reduce the bits-per-channel from 8 bits to 1 bit, the resulting grayscale space is $1/128$ large as the original; the resulting RGB space is only $2^{-(24-3)} = 1/2,097,152$ of the original size. Nevertheless, in Section IV-B we find that squeezing to 4 bits is strong enough to mitigate a lot of adversarial examples while preserving the accuracy on legitimate examples.

2) Implementation

We implement the bit depth reduction operation in Python with the NumPy library. The input and output are in the same numerical scale $[0, 1]$ so that we don't need to change anything of the target models. For reducing to i -bit depth ($1 \leq i \leq 7$), we first multiply the input value with $2^i - 1$ (minus 1 due to the zero value) then round to integers. Next we scale the integers back to $[0, 1]$, divided by $2^i - 1$. The information capacity of the representation is reduced from 8-bit to i -bit with the integer-rounding operation.

B. Spatial Smoothing

Spatial smoothing (also known as *blur*) is a group of techniques widely used in image processing for reducing image noise. Next, we describe the two types of spatial smoothing methods we used: *local smoothing* and *non-local smoothing*.

1) Local Smoothing

Local smoothing methods make use of the nearby pixels to smooth each pixel. By selecting different mechanisms in weighting the neighbouring pixels, a local smoothing method can be designed as Gaussian smoothing, mean smoothing or the median smoothing method [42] we use. As we report in Section IV-C, median smoothing (also known as *median blur* or *median filter*) is particularly effective in mitigating adversarial examples generated by L_0 attacks.

The median filter runs a sliding window over each pixel of the image, where the center pixel is replaced by the median value of the neighboring pixels within the window. It does not actually reduce the number of pixels in the image, but spreads pixel values across nearby pixels. The median filter is essentially squeezing features out of the sample by making adjacent pixels more similar.

The size of the window is a configurable parameter, ranging from 1 up to the image size. If it were set to the image size, it would (modulo edge effects) flatten the entire image to one color. A square shape window is often used in median filtering, though there are other design choices. Several padding methods can be employed for the pixels on the edge, since there are no real pixels to fill the window. We choose

reflect padding [36], in which we mirror the image along with the edge for calculating the median value of a window when necessary.

Median smoothing is particularly effective at removing sparsely-occurring black and white pixels in an image (descriptively known as *salt-and-pepper noise*), whilst preserving edges of objects well.

Figure 4a presents some examples from CIFAR-10 with median smoothing of a 2×2 window in the third column. It suggests why local smoothing can effectively mitigate adversarial examples generated by the Jacobian-based saliency map approach (JSMA) [32] (Section II-B4). JSMA identifies the most influential pixels and modifies their values to a maximum or minimum. The top left is a seed image of the class airplane from the CIFAR-10 dataset. The third image in the first row displays the result of applying a 2×2 median filter to that image. The last row shows the generated adversarial example using the targeted JSMA attack in the leftmost, and the third image illustrates the result of local smoothing of that adversarial example. As with CIFAR-10, both humans and machines see the correct image class clearly after smoothing. We observe the similar effect on the ImageNet dataset with another L_0 attack: CW_0 attack in Figure 4b, even though there are less perturbed pixels.

Implementation. We use the median filter implemented in SciPy [37]. In a 2×2 sliding window, the center pixel is always located in the lower right. When there are two equal-median values due to the even number of pixels in a window, we (arbitrarily) use the greater value as the median.

2) Non-local Smoothing

Non-local smoothing is different from local smoothing because it smooths over similar pixels in a much larger area instead of just nearby pixels. For a given image patch, non-local smoothing finds several similar patches in a large area of the image and replaces the center patch with the average of those similar patches. Assuming that the mean of the noise is zero, averaging the similar patches will cancel out the noise while preserving the edges of an object. Similar with local smoothing, there are several possible ways to weigh the similar patches in the averaging operation, such as Gaussian, mean, and median. We use a variant of the Gaussian kernel because it is widely used and allows to control the deviation from the mean. The parameters of a non-local smoothing method typically include the search window size (a large area for searching similar patches), the patch size and the filter strength (bandwidth of the Gaussian kernel). We will denote a filter as “non-local means (a-b-c)” where “a” means the search window $a \times a$, “b” means the patch size $b \times b$ and “c” means the filter strength.

Figure 4 presents some examples with non-local means (11-3-4). From the first column in Figure 4a, we observe that the adversarial attacks introduce different patterns in the sky background. Non-local smoothing (fourth column) is very effective in restoring the smooth sky while preserving the shape of the airplane. We observe the similar effect from the ImageNet examples in Figure 4b.

Implementation. We use the fast non-local means denoising

method implemented in OpenCV. It first converts a color image to the CIELAB colorspace, then separately denoises its L and AB components, then converts back to the RGB space.

C. Other Squeezing Methods

Our results in this paper are limited to these simple squeezing methods, which are surprisingly effective on our test datasets. However, we believe many other squeezing methods are possible, and continued experimentation will be worthwhile to find the most effective squeezing methods.

One possible area to explore includes lossy compression techniques. Kurakin et al. explored the effectiveness of the JPEG format in mitigating the adversarial examples [20]. Their experiment shows that a very low JPEG quality (e.g. 10 out of 100) is able to destruct the adversarial perturbations generated by FGSM with $\epsilon=16$ (at scale of [0,255]) for at most 30% of the successful adversarial examples. However, they didn’t evaluate the potential loss on the accuracy of legitimate inputs.

Another possible direction is dimension reduction. For example, Turk and Pentland’s early work pointed out that many pixels are irrelevant features in the face recognition tasks, and the face images can be projected to a feature space named *eigenfaces* [40]. Even though image samples represented in the *eigenface*-space lose the spatial information a CNN model needs, the image restoration through *eigenfaces* may be a useful technique to mitigate adversarial perturbations in a face recognition task.

IV. ROBUSTNESS

The previous section demonstrated that images, as used in classification tasks, contain many irrelevant features that can be squeezed without reducing recognizability. For feature squeezing to be effective in detecting adversarial examples (Figure 1), it must satisfy two properties: (1) on adversarial examples, the squeezing reverses the effects of the adversarial perturbations; and (2) on normal legitimate examples, the squeezing does not significantly impact a classifier’s prediction. This section evaluates the how well different feature squeezing methods achieve these properties against various adversarial attacks.

Threat model. In evaluating robustness, we assume a powerful adversary who has full access to a target trained model, but no ability to influence that model. The adversary is not aware of feature squeezing being performed on the operator’s side. With the goal to find inputs that are misclassified by the model, the adversary tries to fool the target model with the white-box attack techniques, whereas the adversarial examples will be inferred by the model with feature squeezing.

We do not propose using feature squeezing directly as a defense because an adversary may take advantage of feature squeezing in attacking a DNN model. For example, when facing binary squeezing, an adversary can construct an image by setting all pixel intensity values to be near 0.5. This image is entirely gray to human eyes. By setting pixel values to either 0.499 or 0.501 it can result in an arbitrary 1-bit filtered image after squeezing, either entirely white or black. Such an attack can easily be detected by our detection framework (Section V), because since the prediction difference between the original and the squeezed will clearly exceed a normal threshold. In

TABLE I: Summary of the target DNN models.

Dataset	Model	Top-1 Accuracy	Top-1 Mean Confidence	Top-5 Accuracy
MNIST	7-Layer CNN [3]	99.43%	99.39%	-
CIFAR-10	DenseNet [18], [23]	94.84%	92.15%	-
ImageNet	MobileNet [17], [24]	68.36%	75.48%	88.25%

more details, we consider how adversaries can adapt to our detection framework in Section V-D.

A. Experimental Setup

We evaluate our defense on state-of-the-art models for the three image datasets, against eleven attack variations representing the best known attacks to date.

Target Models. We use three popular datasets for the image classification task: MNIST, CIFAR-10, and ImageNet. For each dataset, we set up a pre-trained model with the state-of-the-art performance. Table I summarizes the prediction performance of each model and the information of its DNN architecture. Our MNIST model (a seven-layer CNN [3]) achieves a test accuracy of 99.43%; our CIFAR-10 model (a DenseNet [18], [23]) achieves 94.84% test accuracy. The prediction performance of both models is competitive with state-of-the-art results [1]. For the ImageNet dataset, we use a MobileNet model [17], [24] because MobileNet is more widely used on mobile phones and its small and efficient design make it easier to conduct experiments. The pre-trained MobileNet model achieves top-1 accuracy 68.36% and top-5 accuracy 88.25%, both are comparable to state-of-the-art results. In contrast, a larger model such as Inception v3 [38], [7] with six times of trainable parameters could achieve top-1 accuracy 76.28% and top-5 accuracy 93.03%. However, the calculation on such a model is much more expensive due to the massive architecture.

Attacks. We evaluate feature squeezing on all of the attacks described in Section II-C. For the targeted attacks, we try each attack with two types of targets: the next class (-Next),

$$t = L + 1 \pmod{\#classes}, \quad (7)$$

and the least-likely class (-LL),

$$t = \min(\hat{\mathbf{y}}), \quad (8)$$

Here t is the target class, L is the index of the ground-truth class and $\hat{\mathbf{y}}$ is the prediction vector of an input image. This gives eleven total attacks: the three untargeted attacks (FGSM, BIM and DeepFool), and two versions each of the four targeted attacks (JSMA, CW_∞ , CW_2 , and CW_0). We use the implementations of FGSM, BIM and JSMA provided by the Cleverhans library [30]. For DeepFool and the three CW attacks, we use the implementations from the original authors [3], [27]. The parameters we use for the attacks are given in Table VI (in the appendix).¹

For the seed images, we select the first 100 correctly predicted examples in the test (or validation) set from each dataset for all the attack methods, since some attacks are too expensive

¹All of our models and codes for attacks, defenses, and testing are available as an open source tool (<https://github.com/mzweilin/EvadeML-Zoo>).

TABLE II: Evaluation of 11 different attacks (each with 100 seed images) against DNN models on three datasets. The *cost* of an attack generating adversarial examples is measured in seconds per sample. The L_0 distortion is normalized by the number of pixels (e.g., 0.56 means 56% of all pixels in the image are modified).

	Configuration		Cost	Success Rate	Prediction Confidence	Distortion				
	Attack	Mode				L_∞	L_2	L_0		
MNIST	L_∞	FGSM	0.002	46%	93.89%	0.3020	5.9047	.5601		
		BIM	0.01	91%	99.62%	0.3020	4.7580	.5132		
		CW_∞	Next	51.25	100%	99.99%	0.2513	4.0911	.4906	
			LL	49.95	100%	99.98%	0.2778	4.6203	.5063	
		L_2	CW_2	Next	0.33	99%	99.23%	0.6556	2.8664	.4398
				LL	0.38	100%	99.99%	0.7342	3.2176	.4362
	L_0	CW_0	Next	68.76	100%	99.99%	0.9964	4.5378	.0473	
			LL	74.55	100%	99.99%	0.9964	5.1064	.0597	
		JSMA	Next	0.79	71%	74.52%	1.0000	4.3276	.0473	
			LL	0.98	48%	74.80%	1.0000	4.5649	.0535	
	CIFAR-10	L_∞	FGSM	0.02	85%	84.85%	0.0157	0.8626	.9974	
			BIM	0.19	92%	95.29%	0.0078	0.3682	.9932	
CW_∞			Next	225.32	100%	98.22%	0.0122	0.4462	.9896	
			LL	224.58	100%	97.79%	0.0143	0.5269	.9947	
L_2		DeepFool	0.36	98%	73.45%	0.0279	0.2346	.9952		
		CW_2	Next	10.36	100%	97.90%	0.0340	0.2881	.7677	
			LL	12.01	100%	97.35%	0.0416	0.3577	.8549	
		L_0	CW_0	Next	366.54	100%	98.19%	0.6500	2.1033	.0186
LL				426.05	100%	97.60%	0.7121	2.5300	.0241	
JSMA			Next	8.44	100%	43.29%	0.8960	4.9543	.0790	
			LL	13.64	98%	39.75%	0.9037	5.4883	.0983	
ImageNet		L_∞	FGSM	0.02	99%	63.99%	0.0078	3.0089	.9941	
	BIM		0.18	100%	99.71%	0.0039	1.4059	.9839		
	CW_∞		Next	210.70	99%	90.33%	0.0059	1.3118	.8502	
			LL	268.86	99%	81.42%	0.0095	1.9089	.9520	
	L_2	DeepFool	60.16	89%	79.59%	0.0269	0.7258	.9839		
		CW_2	Next	20.63	90%	76.25%	0.0195	0.6663	.3226	
			LL	29.14	97%	76.03%	0.0310	1.0267	.5426	
		L_0	CW_0	Next	607.94	100%	91.78%	0.8985	6.8254	.0030
	LL			979.05	100%	80.67%	0.9200	9.0816	.0053	

to run on all the seeds. We adjust the applicable parameters of each attack to generate high-confidence adversarial examples, otherwise they would be easily rejected. This is because the three DNN models we use achieve high confidence of the top-1 predictions on legitimate examples (see Table I; mean confidence is over 99% for MNIST, 92% for CIFAR-10, and 75% for ImageNet). In addition, all the pixel values in the generated adversarial images are clipped and squeezed to 8-bit-per-channel pixels so that the resulting inputs are within the possible space of images.

We use a PC equipped with an i7-6850K 3.60GHz CPU and 64GiB system memory as well as a GeForce GTX 1080 to conduct the experiments.

In Table II, we evaluate the adversarial examples regarding the success rate, the run-time cost, the prediction confidence and the distance to the seed image measured by L_2 , L_∞ and L_0 metrics. The evaluation results for all eleven attacks on the three datasets are provided. The success rate captures the probability an adversary achieves their goal. For untargeted attacks, the success rate is calculated as $1 - accuracy$; for targeted attacks, it is the accuracy for the targeted class. Table II shows that in general most attacks generate high-confidence adversarial examples against three DNN models with a high success rate. The CW attacks often produce fewer distortions than other attacks using the same norm objective but are much more expensive to generate. On the other hand, FGSM, DeepFool, and JSMA often produce low-confidence adversarial examples. We exclude the DeepFool attack from the MNIST dataset because it generates images that appear

unrecognizable to human eyes. We do not have JSMA results for the ImageNet dataset because the available implementation ran out of memory on our 64GiB test machine.

In Table III we evaluate and compare how different feature squeezers influence the classification accuracy of DNN models on three image datasets for all attacks. We discuss experimental results of each type of squeezers further below.

B. Color Depth Reduction

The resolution of a specific bit depth is defined as the number of possible values for each pixel. For example, the resolution of 8-bit color depth is 256. Reducing the bit depth lowers the resolution and diminishes the opportunity an adversary has to find effective perturbations. Since an adversary’s goal is to produce small and imperceptible perturbations in the case of adversarial examples, as the resolution is reduced, such small perturbations no longer have any impact.

MNIST. The Last column of Table III shows the binary filter (1-bit depth reduction) barely reduces the accuracy on the legitimate examples of MNIST (from 99.43% to 99.33% on the test set). When comparing the model accuracy on the adversarial examples by the original classifier (the first row with squeezer None) to the one with the binary filter (the second row with squeezer bit depth (1-bit)), we see the binary filter is effective on all the L_2 and L_∞ attacks. For example, it improves the accuracy on CW_∞ adversarial examples from 0% to 100%. Interestingly the binary filter works well even for large L_∞ distortions. This is because the binary filter squeezes each pixel into 0 or 1 using a cutoff 0.5 in the $[0, 1)$ scale. This means maliciously perturbing a pixel’s value by ± 0.30 has no affect on those pixels whose original values fall into $[0, .20)$ and $[.80, 1)$. In contrast, bit depth reduction is not effective against L_0 attacks (JSMA and CW_0) since these attacks make large changes to a few pixels and can not be reversed by the bit depth squeezer. The next section shows that the spatial smoothing squeezers are often effective against L_0 attacks.

CIFAR-10 and ImageNet. Because the DNN models for CIFAR-10 and ImageNet are more sensitive to the adversary, adversarial examples at very low L_2 and L_∞ distortions can be found. Table III includes the results of 4-bit depth and 5-bit depth filters in mitigating the adversaries for CIFAR-10 and ImageNet. The 5-bit depth in testing increases the accuracy on adversarial inputs for several of the attacks (for example, increasing accuracy from 0% to 40% for the CW_2 next class targeted attack), while almost perfectly preserving the accuracy on legitimate data (94.55% compared with 94.84%). The more aggressive 4-bit depth filter is more robust against adversaries. For example, the accuracy on CW_2 increases to 84%, but it reduces the accuracy on legitimate inputs from 94.84% to 93.11%. We do not believe these results are good enough for use as a stand-alone defense (even ignoring the risk of adversarial adaptation), but they provide some insight why the method is effective as used in our detection framework.

C. Median Smoothing

The adversarial perturbations produced by the L_0 attacks (JSMA and CW_0) are similar to *salt-and-pepper noise*, though it is introduced intentionally instead of randomly. Note that

the adversarial strength of an L_0 adversary limits the number of pixels that can be manipulated, so it is not surprising that maximizing the amount of change to each modified pixel is typically most useful to the adversary. This is why the smoothing squeezers are more effective against these attacks than the color depth squeezers.

MNIST. We evaluate two window sizes on the MNIST dataset in Table III. Median smoothing is the best squeezer for all of the L_0 attacks (CW_0 and JSMA). The median filter with 2×2 window size performs slightly worse on adversarial examples than the one with 3×3 window, but it almost perfectly preserves the performance on the legitimate examples (decreasing accuracy from 99.43% to 99.28%).

CIFAR-10 and ImageNet. The experiment confirms the intuition suggested by Figure 4a that median smoothing can effectively eliminate the L_0 -limited perturbations. Without squeezing, the L_0 attacks are effective on CIFAR-10, resulting in 0% accuracy for the original model (“None” row in Table III). However, with a 2×2 median filter, the accuracy increases to over 75% for all the four L_0 type attacks. We observe similar results on ImageNet, where the accuracy increases from 0% to 85% for the CW_0 attacks after median smoothing.

D. Non-local Smoothing

The image examples in Figure 4a suggest that non-local smoothing is inferior to median smoothing in eliminating the L_0 type perturbations, but superior for smoothing the background and preserving the object edges. This intuition is confirmed by the experimental results on CIFAR-10 and ImageNet (because the MNIST images are hand-drawn digits that are not conducive to finding similar patches, we do not consider non-local smoothing on MNIST). From Table III we learn that non-local smoothing has comparable performance in increasing the accuracy on adversarial examples other than the L_0 type. On the other hand, it has little impact on the accuracy on legitimate examples. For example, the 2×2 median filter decreases the accuracy on the CIFAR-10 model from 94.84% to 89.29% while the model with non-local smoothing still achieves 91.18%. We do not apply the non-local smoothing on MNIST images because it is difficult to find similar patches on such images for smoothing a center patch.

E. Combining with Adversarial Training

Since our approach modifies inputs rather than the model, it is compatible with any defense technique that operates on the model. The most successful previous defense against adversarial examples is adversarial training (Section II-C). To evaluate the effectiveness of composing our feature squeezing method with adversarial training, we combined it with the adversarial training implemented by Cleverhans [30]. The objective is to minimize the mean loss on the legitimate examples and the adversarial ones generated by FGSM on the fly with $\epsilon = 0.3$. The model is trained in 100 epochs.

Figure 5 shows that the bit depth reduction by itself significantly outperforms the adversarial training method on MNIST in face of the FGSM adversary, but that composing both methods produces even better results. Used by itself, the binary filter feature squeezing outperforms adversarial training

TABLE III: Model accuracy with feature squeezing

Dataset	Squeezer		L_∞ Attacks				L_2 Attacks				L_0 Attacks				All Attacks	Legitimate
	Name	Parameters	FGSM	BIM	CW $_\infty$		Deep-Fool	CW $_2$		CW $_0$		JSMA				
					Next	LL		Next	LL	Next	LL	Next	LL			
MNIST	None		54%	9%	0%	0%	-	0%	0%	0%	0%	27%	40%	13.00%	99.43%	
	Bit Depth	1-bit	92%	87%	100%	100%	-	83%	66%	0%	0%	50%	49%	62.70%	99.33%	
	Median Smoothing	2x2	61%	16%	70%	55%	-	51%	35%	39%	36%	62%	56%	48.10%	99.28%	
		3x3	59%	14%	43%	46%	-	51%	53%	67%	59%	82%	79%	55.30%	98.95%	
CIFAR-10	None		15%	8%	0%	0%	2%	0%	0%	0%	0%	0%	0%	2.27%	94.84%	
	Bit Depth	5-bit	17%	13%	12%	19%	40%	40%	47%	0%	0%	21%	17%	20.55%	94.55%	
		4-bit	21%	29%	69%	74%	72%	84%	84%	7%	10%	23%	20%	44.82%	93.11%	
	Median Smoothing	2x2	38%	56%	84%	86%	83%	87%	83%	88%	85%	84%	76%	77.27%	89.29%	
	Non-local Means	11-3-4	27%	46%	80%	84%	76%	84%	88%	11%	11%	44%	32%	53.00%	91.18%	
ImageNet	None		1%	0%	0%	0%	11%	10%	3%	0%	0%	-	-	2.78%	69.70%	
	Bit Depth	4-bit	5%	4%	66%	79%	44%	84%	82%	38%	67%	-	-	52.11%	68.00%	
		5-bit	2%	0%	33%	60%	21%	68%	66%	7%	18%	-	-	30.56%	69.40%	
	Median Smoothing	2x2	22%	28%	75%	81%	72%	81%	84%	85%	85%	-	-	68.11%	65.40%	
		3x3	33%	41%	73%	76%	66%	77%	79%	81%	79%	-	-	67.22%	62.10%	
Non-local Means	11-3-4	10%	25%	77%	82%	57%	87%	86%	43%	47%	-	-	57.11%	65.40%		

No results are shown for DeepFool on MNIST because of the adversarial examples it generates appear unrecognizable to humans; no results are shown for JSMA on ImageNet because it requires more memory than available to run.

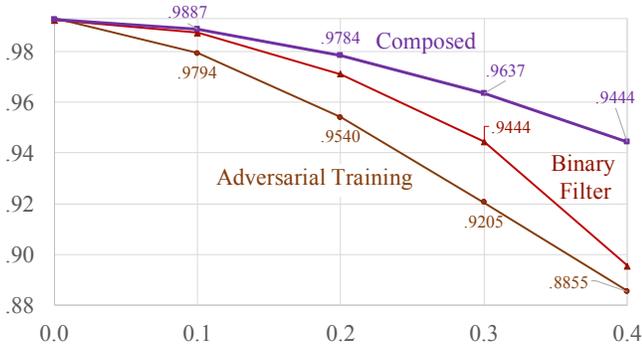


Fig. 5: Composing adversarial training with feature squeezing. The horizontal axis is ϵ , so the adversarial strength increases to the right. By itself, bit depth reduction on the original model outperforms adversarial training. The adversarial-trained model is fed with the original training examples as well as those generated by FGSM ($\epsilon = 0.3$) during the 100-epoch training phase. Composing the 1-bit filter with the adversarial-trained model performs even better.

for ϵ values ranging from 0.1 to 0.4². This is the best case for adversarial training since the adversarially-trained model is learning from the same exact adversarial method (retraining is done with FGSM examples generated at $\epsilon = 0.3$) as the one used to produce the adversarial examples in the test. Nevertheless, feature squeezing still outperforms it, even at the same $\epsilon = 0.3$ value: 94.44% accuracy on adversarial examples compared to 92.05%.

Feature squeezing is far less expensive than adversarial training. It is almost cost-free, as we simply insert a binary filter before the pre-trained MNIST model. On the other hand, adversarial training is very expensive as it requires both generating adversarial examples and retraining the classifier for many epochs.³

²The choice of ϵ is arbitrary, but examples where $\epsilon > 0.3$ are typically not considered valid adversarial examples [10] since such high ϵ values produce images that are obviously different from the original images.

³We would like to test retraining with the stronger adversaries, and on the CIFAR-10 and ImageNet datasets also, but have not been able to do this experiment as the time to do adversarial training on larger models is prohibitively expensive.

When its cost is not prohibitive, though, adversarial training is still beneficial since it can be combined with feature squeezing. Simply inserting a binary filter before the adversarially-trained model increases the robustness against an FGSM adversary. Figure 5 shows that the accuracy on adversarial inputs with $\epsilon = 0.3$ is 96.37% for the combined model, which significantly outperforms both standalone approaches: 92.05% for adversarial training and 94.44% for the bit depth reduction.

V. DETECTING ADVERSARIAL INPUTS

From Section IV we see that feature squeezing is capable of obtaining accurate model predictions for many adversarial examples with little reduction in accuracy for legitimate examples. This enables detection of adversarial inputs using the framework introduced in Figure 1. The basic idea is to compare the model’s prediction on the original sample with the same model’s prediction on the sample after squeezing. The model’s predictions for a legitimate example and its squeezed version should be similar. On the contrary, if the original and squeezed examples result in dramatically different predictions, the input is likely to be adversarial. Table IV and Figure 6 summarize the results of our experiments that confirm this intuition for all three datasets. The following subsections provide more details on our detection method, experimental setup, and discuss the results. Section V-D considers how adversaries may adapt to our defense.

A. Detection Method

A prediction vector generated by a DNN classifier normally represents the probability distribution how likely an input sample is to belong to each possible class. Hence, comparing the model’s original prediction with the prediction on the squeezed sample involves comparing two probability distribution vectors. There exist several ways to compare the probability distributions, such as the L_1 norm, the L_2 norm and K-L divergence [2]. For this work, we select the L_1 norm⁴ as a natural measure of the difference between the original

⁴This turned out to work well, but it is certainly worth exploring in future work if other metrics can work better.

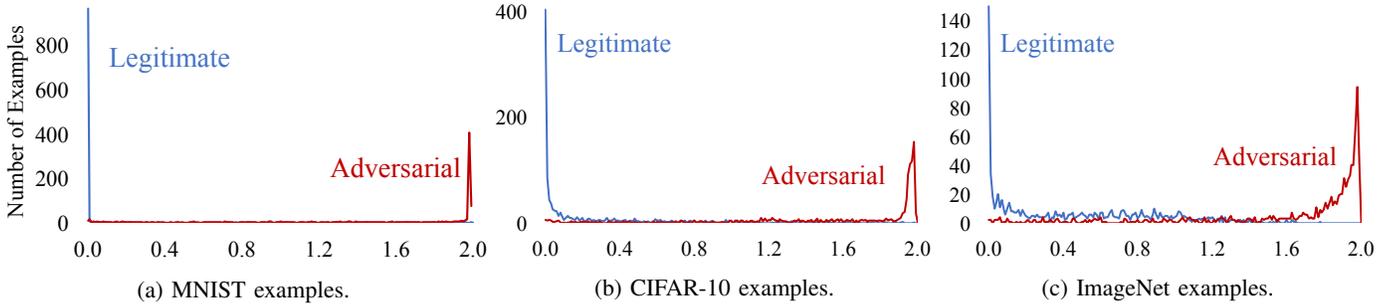


Fig. 6: Differences in L_1 distance between original and squeezed sample, for legitimate and adversarial examples across three datasets. The L_1 score has a range from 0.0 to 2.0. Each curve is fitted over 200 histogram bins each representing the L_1 distance range of 0.01. Each sample is counted in the bin for the maximum L_1 distance between the original prediction and the output of the best joint-detection squeezing configuration shown in Table IV. The curves for adversarial examples are for all adversarial examples, including unsuccessful ones (so the separation for successful ones is even larger than shown here).

prediction vector and the squeezed prediction:

$$score^{(\mathbf{x}, \mathbf{x}_{squeezed})} = \|g(\mathbf{x}) - g(\mathbf{x}_{squeezed})\|_1 \quad (9)$$

Here $g(\mathbf{x})$ is the output vector of a DNN model produced by the softmax layer whose i^{th} entry describes the probability how likely input \mathbf{x} is in the i -th class.

The L_1 norm function is defined as $\|z\|_1 = \sum_i |z_i|$. The L_1 score has a range from 0 to 2 for the prediction vectors. A higher score means there is a greater difference between the original prediction and the squeezed prediction for an input \mathbf{x} . The maximum value of 2 is reached when each prediction vector consists of a 1 and all zeros, but with different classes as the 1. Based on the accuracy results in Section IV, we expect the score to be small for legitimate inputs and large for adversarial examples. The effectiveness of detection depends on selecting a threshold value that accurately distinguishes between legitimate and adversarial inputs.

Even though we can select an effective feature squeezer for a specific type of adversarial method, an operator typically does not know the exact attack method that would be used in practice. Hence, we combine multiple feature squeezers for detection by outputting the maximum distance:

$$score^{joint} = \max(score^{(\mathbf{x}, \mathbf{x}_{sq1}), score^{(\mathbf{x}, \mathbf{x}_{sq2}), \dots}) \quad (10)$$

We choose the max operator based on an assumption that there will be a squeezer being effective on one adversarial example which generates the highest L_1 score. On the other hand, it may increase the false positive rate because the max operator also selects the most destructive squeezer on legitimate inputs. We observed that we could usually make a reasonable trade-off in the empirical results. We may investigate a better way to combine the squeezers in the future work.

Figure 6 shows the histogram of $score^{joint}$ for both legitimate (blue) and adversarial examples (red) on three datasets. The peak for legitimate examples is always near 0, and the peak for adversarial examples is always near 2. Picking a threshold value between the two peaks is a balance between high detection rates and acceptable false positive rates. For our experiments, we require a false positive rate below 5% (see Section V-B).

B. Experimental Setup

We report on experiments using all attacks from Section IV with the three types of squeezers in different configurations.

Datasets. To get a balanced dataset for detection, we select the same number of legitimate examples from the test (or validation) set of each dataset. For each of the attacks in Section IV, we use the 100 adversarial examples generated for each attack in the robustness experiments. This results in 2,000 total examples for MNIST (of which 1,000 are legitimate examples, and 1,000 are adversarial), 2,200 examples for CIFAR-10 and 1,800 examples for ImageNet. We randomly split each detection dataset into two groups: one-half for training the detector and the remainder for validation. Note that some of the adversarial examples are failed adversarial examples, that do not confuse the original model, so the number of successful adversarial examples varies slightly across the attacks.

Squeezers. We first consider the artificial situation where the defender knows the attack method, and evaluate how well each squeezing configuration does against adversarial examples generated by each attack method. Then, we consider the realistic scenario where the defender does not know that attack method used by the adversary and needs to select a configuration that works well against a distribution of possible attacks.

Training. The training phase of our detector is simply selecting an optimal threshold of $score^{joint}$. One typical practice is to find the one that maximizes the training accuracy. However, a detector with high accuracy could be useless in many security-sensitive tasks if it had a high false positive rate since the actual distribution of samples is not balanced and mostly benign. Therefore, we instead select a threshold that ensures the false positive rate below 5%, choosing a threshold that is exceeded by just below 5% of legitimate samples. Note that the training threshold is set using only the legitimate examples, so does not depend on the adversarial examples.

Validation. Next, we use the chosen threshold value to measure the detection rate on three groups: the successful adversarial examples (SAEs), the failed adversarial examples (FAEs), and the legitimate examples (for false positive rate).

Except when noted explicitly, “detection rate” means the detection rate on successful adversarial examples. We think it is important to distinguish failed adversarial examples from legitimate examples here since detecting failed adversarial examples is useful for detecting attacks early, whereas an alarm on a legitimate example is always undesirable and is counted as a false positive.

C. Results

First, we discuss the effectiveness of different squeezers against different attacks. Then, we consider how multiple squeezers would work in a joint detection framework.

Detection with Single Squeezer. Table IV summarizes the validation results for detection on each dataset, showing the detection rates for successful adversarial examples for each attack method with a variety of configurations. For each dataset, we first list the detection rate of several detectors built upon single squeezers. For each squeezing method, we tried several parameters and compare the performance for each dataset. The “Best Attack-Specific Single Squeezer” row gives the detection rate for the best single squeezer against a particular attack. This represents the (unrealistically optimistic) case where the model operator knows the attack type and selects a single squeezer for detection that may be different for each attack. Below this, we show the best result of joint detection (to be discussed later) with multiple squeezers where the same configuration is used for all attacks.

The best bit depth reduction for MNIST is squeezing the color bits to one, which achieves 100% detection rate for 6 of the attacks and at least 98.9% detection for all the L_∞ and L_2 attacks. It is not as effective on the advanced CW_0 attacks, however, since these attacks are making large changes to a small number of pixels. On the contrary, the 3×3 median smoothing is the most effective on detecting the L_0 attacks with detection rates above 92.00%. This matches the observation from Table III that they have different strengths for improving the model accuracy. For MNIST, there is at least one squeezer that provides good (above 92% detection) results for all of the attacks.

For CIFAR-10, we find that 2×2 median smoothing is the best single squeezer for detecting every attack except DeepFool, which is best detected by non-local smoothing. This is consistent with the robustness results in Table III. For the ImageNet dataset, we find several different squeezers are similarly effective on each attack type. For example, the CW_2 -LL attack can be 100% detected with several bit depth filters, the 2×2 median smoothing or some non-local mean filters.

The third column in the table gives the distance threshold setting needed to satisfy the maximum false positive rate of 5%. These threshold values provide some insight into how well a particular squeezer distinguishes adversarial from legitimate examples. For the binary filter on MNIST, a tiny threshold value of 0.0005 was sufficient to produce a false positive rate below 5%, which means the squeezing has negligible impact on the legitimate examples: 95% of the legitimate examples have the L_1 -based distance score below 0.0005. On the other hand, the best median smoothing filter (2×2) on MNIST needs a larger threshold value 0.0029 to achieve a similar false

positive rate, which means it is slightly more destructive than the binary filter on the legitimate examples. The more aggressive median smoothing with 3×3 window results in an even higher threshold 0.0390, because the legitimate examples could get over-squeezed to the target classifier. A lower threshold is always preferred for detection, which makes it more sensitive to adversarial examples.

For some of the attacks, none of the feature squeezing methods work well enough for the color datasets. The worst cases, surprisingly, are for FGSM and BIM, two of the earlier adversarial methods. The best single-squeezer-detection only recognizes 25.88% of the successful FGSM examples and 52.17% of BIM on the CIFAR-10 dataset, while the detection rates are 44.41% and 59.00% on ImageNet. We suspect the reason the tested squeezers are less effective against these attacks is because they make larger perturbations than the more advanced attacks (especially the CW attacks), and the feature squeezers we use are well suited to mitigating small perturbations. Understanding why these detection rates are so much lower than the others, and developing feature squeezing methods that work well against these attacks is an important avenue for future research.

Joint-Detection with Multiple Squeezers. Table V summarizes the overall detection rates for the best joint detectors against a distribution of all adversarial methods. By comparing the last two rows of each dataset in Table IV, we see that joint-detection often outperforms the best detector with a single squeezer. For example, the best single-squeezer-detection detects 99.00% of the CW_2 -LL examples while the joint detection makes it 100%.

The main reason to use multiple squeezers, however, is because this is necessary to detect unknown attacks. Since the model operator is unlikely to know what attack adversaries may use, it is important to be able to set up the detection system to work well against any attack. For each data set, we try several combinations of the three squeezers with different parameters and find out the configuration that has the best detection results across all the adversarial methods (shown as the “Best Joint Detection” in Table IV). For MNIST, the best combination was the 1-bit depth squeezer with 2×2 median smoothing, combining the best parameters for each type of squeezer. For the color image datasets, different combinations were found to outperform combining the best squeezers of each type. For example, the best joint detection configuration for ImageNet includes the 5-bit depth squeezer, even though the 3-bit depth squeezer was better as a single squeezer.

Comparing the “Best Attack-Specific Single Squeezer” and “Best Joint Detection” rows in Table IV reveals that the joint detection is usually competitive with the best single squeezers over all the attacks. Since the joint detector needs to maintain the 5% false positive rate requirement, it has a higher threshold than the individual squeezers. This means in some cases its detection rate for a particular attack will be worse than the best single squeezer achieves. For MNIST, the biggest drop is for detection rate for CW_0 (LL) attacks drops from 98% to 92%; for CIFAR-10, the joint squeezer always outperforms the best single squeezer; for ImageNet, the detection rate drops for BIM (59% to 52%), CW_2 (Next) (93% to 90%), and CW_0 (Next) (99% to 98%). For simplicity, we use a single threshold

TABLE IV: Detection rate on successful adversarial examples.

	Configuration			L_∞ Attacks				L_2 Attacks			L_0 Attacks				Overall Detection Rate
	Squeezer	Parameters	Threshold	FGSM	BIM	CW_∞		Deep Fool	CW_2		CW_0		JSMA		
						Next	LL		Next	LL	Next	LL	Next	LL	
MNIST	Bit Depth	1-bit*	0.0005	100.00%	98.90%	100.00%	100.00%	-	100.00%	99.00%	54.00%	57.00%	100.00%	100.00%	90.30%
		2-bit	0.0002	67.39%	7.69%	62.00%	76.00%	-	93.94%	90.00%	40.00%	41.00%	98.59%	100.00%	65.59%
	Median Smoothing	2x2*	0.0029	73.91%	26.37%	100.00%	100.00%	-	96.97%	99.00%	84.00%	91.00%	97.18%	100.00%	86.84%
		3x3	0.0390	47.83%	12.09%	80.00%	83.00%	-	83.84%	93.00%	92.00%	98.00%	98.59%	100.00%	78.06%
	Best Attack-Specific Single Squeezer		-	100.00%	98.90%	100.00%	100.00%	-	100.00%	99.00%	92.00%	98.00%	100.00%	100.00%	-
Best Joint Detection (1-bit, 2x2)		0.0029	100.00%	98.90%	100.00%	100.00%	-	100.00%	100.00%	93.00%	92.00%	100.00%	100.00%	98.15%	
CIFAR-10	Bit Depth	1-bit	1.9997	4.71%	8.70%	0.00%	0.00%	1.02%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	1.29%
		2-bit	1.9967	7.06%	15.22%	1.00%	0.00%	0.00%	1.00%	0.00%	0.00%	1.00%	0.00%	0.00%	2.22%
		3-bit	1.7822	9.41%	30.43%	78.00%	97.00%	19.39%	85.00%	96.00%	31.00%	30.00%	0.00%	0.00%	41.04%
		4-bit*	0.7930	11.76%	26.09%	83.00%	90.00%	65.31%	95.00%	98.00%	12.00%	21.00%	3.00%	4.08%	45.29%
		5-bit	0.3301	3.53%	11.96%	42.00%	63.00%	55.10%	76.00%	86.00%	5.00%	3.00%	7.00%	6.12%	31.42%
	Median Smoothing	2x2*	1.1296	25.88%	52.17%	96.00%	100.00%	71.43%	98.00%	100.00%	98.00%	100.00%	81.00%	85.71%	84.29%
	Non-local Means	3x3	1.9431	5.88%	23.91%	70.00%	90.00%	2.04%	69.00%	89.00%	74.00%	95.00%	4.00%	9.18%	48.80%
		11-3-2	0.2777	17.65%	36.96%	85.00%	93.00%	79.59%	91.00%	96.00%	8.00%	7.00%	30.00%	21.43%	48.80%
		11-3-4	0.7537	21.18%	48.91%	91.00%	96.00%	74.49%	95.00%	100.00%	25.00%	35.00%	29.00%	25.51%	55.45%
		13-3-2	0.2904	17.65%	35.87%	87.00%	94.00%	78.57%	91.00%	96.00%	8.00%	8.00%	32.00%	24.49%	49.17%
		13-3-4*	0.8290	21.18%	47.83%	91.00%	96.00%	72.45%	95.00%	100.00%	26.00%	34.00%	27.00%	24.49%	54.90%
	Best Attack-Specific Single Squeezer		-	25.88%	52.17%	96.00%	100.00%	79.59%	98.00%	100.00%	98.00%	100.00%	81.00%	85.71%	-
Best Joint Detection (5-bit, 2x2, 13-3-2)		1.1402	27.06%	52.17%	97.00%	100.00%	79.59%	99.00%	100.00%	98.00%	100.00%	81.00%	85.71%	85.03%	
ImageNet	Bit Depth	1-bit	1.9957	0.00%	0.00%	1.01%	1.01%	1.12%	0.00%	0.00%	0.00%	0.00%	-	-	6.90%
		2-bit	1.9301	16.16%	53.00%	58.59%	45.45%	33.71%	24.44%	38.14%	38.00%	28.00%	-	-	40.00%
		3-bit*	1.4293	13.13%	48.00%	95.96%	100.00%	62.92%	83.33%	100.00%	85.00%	100.00%	-	-	76.43%
		4-bit	0.7947	9.09%	10.00%	84.85%	100.00%	47.19%	88.89%	100.00%	76.00%	100.00%	-	-	68.57%
		5-bit	0.3686	6.06%	2.00%	64.65%	96.97%	32.58%	92.22%	100.00%	46.00%	99.00%	-	-	59.29%
	Median Smoothing	2x2*	1.0985	41.41%	41.00%	95.96%	100.00%	76.40%	92.22%	100.00%	99.00%	100.00%	-	-	84.05%
	Non-local Means	3x3	1.4348	37.37%	59.00%	92.93%	97.98%	69.66%	83.33%	98.97%	97.00%	100.00%	-	-	83.10%
		11-3-2	0.6483	13.13%	13.00%	79.80%	94.95%	42.70%	92.22%	97.94%	90.00%	90.00%	-	-	62.62%
		11-3-4	1.0423	22.22%	36.00%	92.93%	100.00%	61.80%	93.33%	100.00%	73.00%	98.00%	-	-	75.00%
		13-3-2	0.6895	12.12%	13.00%	81.82%	94.95%	42.70%	90.00%	97.94%	48.00%	90.00%	-	-	62.86%
		13-3-4*	1.0809	22.22%	34.00%	94.95%	100.00%	60.67%	93.33%	100.00%	75.00%	98.00%	-	-	75.24%
	Best Attack-Specific Single Squeezer		-	41.41%	59.00%	95.96%	100.00%	76.40%	93.33%	100.00%	99.00%	100.00%	-	-	-
Best Joint Detection (5-bit, 2x2, 11-3-4)		1.2476	44.44%	52.00%	96.97%	100.00%	79.78%	90.00%	100.00%	98.00%	100.00%	-	-	85.24%	

TABLE V: Details of the best joint detectors. SAE: successful adversarial example. FAE: failed adversarial example.

Dataset	Squeezers	Threshold	Detection Rate (SAEs)	Detection Rate (FAEs)	False Positive Rate	ROC-AUC (including FAEs)	ROC-AUC (excluding FAEs)
MNIST	Bit Depth (1-bit), Median (2 × 2)	0.0029	98.15%	20.00%	3.98%	94.51%	99.61%
CIFAR-10	Bit Depth (5-bit), Median (2 × 2), Non-local Means (13-3-2)	1.1402	85.03%	9.09%	4.93%	95.67%	95.86%
ImageNet	Bit Depth (5-bit), Median (2 × 2), Non-local Means (11-3-4)	1.2476	85.24%	25.00%	4.70%	94.04%	94.51%

across all of the squeezers in a joint detector; we expect there are better ways to combine multiple squeezers that would use different thresholds for each of the squeezers to avoid this detection reduction, and plan to study this in future work.

Table V summarizes the overall performance of the best joint detection configuration we found for each dataset. The overall detection rate is 98.15% on MNIST, 85.03% on CIFAR-10 and 85.24% on ImageNet.

We report ROC-AUC scores in Table V both counting the failed adversarial examples as legitimate (Include FAEs), and excluding the failed adversarial examples from consideration (Exclude FAEs), since it is not clear what the correct output should be for a failed adversarial example and in some scenarios detecting a failed adversarial example as an attack is beneficial. Our joint-detector achieves around 95% ROC-AUC score for all three datasets. Excluding the failed adversarial examples, the ROC-AUC of the detector is as high as 99.61% for MNIST, and above 94.5% for ImageNet. The false positive rate on legitimate examples are all lower than 5%, which is expected considering how we select a threshold value in the training phase. The detection rate for the best configuration on successful adversarial examples exceeds 98% for MNIST using a 1-bit filter and a 2 × 2 median filter and exceeds 85% for the

other two datasets using a combination of three types feature squeezing methods with different parameters. The detection rates for failed adversarial examples are much lower than those for successful adversarial examples, but much higher than the false positive rate for legitimate examples. This is unsurprising since FAEs are attempted adversarial examples, but since they are not successful the prediction outputs for the unsqueezed and squeezed inputs are more similar.

D. Adversarial Adaptation

So far, we have only considered static adversaries, who use state-of-the-art methods to find adversarial examples but do not adapt to attack our feature squeezing method directly. It is difficult to speculate on future attacks, but in this section we consider the challenges an adversary faces in adapting attacks to our defenses and report on preliminary experiments with adaptive attacks.

To be successful against our detection framework, an adversary needs to find an input where the original classifier produces the wrong output and the L_1 score between the model’s predictions on squeezed and original inputs is below the detection threshold. This is a much harder problem than just finding an adversarial example, as is supported by our

experimental results.

He et al. [14] recently proposed an adaptive attack which can successfully find adversarial examples that defeat one configuration of a feature squeezing defense.⁵ The approach finds adversarial examples that both confuse the original model and have a $score^{joint}$ lower than a pre-selected threshold for squeezed inputs. Their approach adapts the CW_2 attack by adding a penalty term for the L_1 prediction distance. It requires that all the feature squeezing operators are differentiable so that it is possible to compute the gradient of the loss function in the optimization process. For the non-differentiable feature squeezers such as the bit depth reduction, their approach requires restarting the algorithm several times with random initialization and hoping it finds an example that is resilient against the non-differentiable squeezers. Due to this reason, the attack is non-deterministic and more time-consuming in face of non-differentiable components [14]. The attack takes roughly 20 seconds on each MNIST sample, which is around $60\times$ slower than the original CW_2 attack.

We only evaluate their adaptive attack on the MNIST dataset, because we currently don't have a Tensorflow implementation of the non-local means filter used on CIFAR-10 and ImageNet.⁶ We adapted their attacking code to the MNIST model in this paper and the joint-detector we report in Table V with the threshold value 0.0029.⁷

In addition to the untargeted attack He et al. tested, we also test the two targeted attacks mentioned earlier in Section IV. We used the same parameters as He et al. [14] in generating the adaptive adversarial examples except that we change the target threshold value to 0.0029 that is used by the detector in this paper. The seed examples are the same 100 images used in Section IV.

With the exact threshold we use in this paper as the target value, the success rate of the untargeted attack is 64%. This is lower than the 100% success rate in their report. The success rates of the targeted attacks are even lower: 41% for the targeted-next attack and 27% for the targeted-ll attack. We believe these lower adversarial success rates are due to using a different MNIST model (representative of a state-of-the-art model) and our detector has a much smaller threshold value (0.0029) compared to the one they used (over 0.3076).

The adaptive attack often produces adversarial examples with large perturbations that look ambiguous even to human eyes. In order to understand the adversary success rate in a more realistic setting, we clipped the perturbations with different ϵ values ranging from 0 to 1 and presented the success rates in Figure 7. We found that the success rates decreased to 15%, 4%, 0% respectively if we clipped the perturbations with 0.3, which was an ϵ value considered in many papers for the MNIST dataset[22].

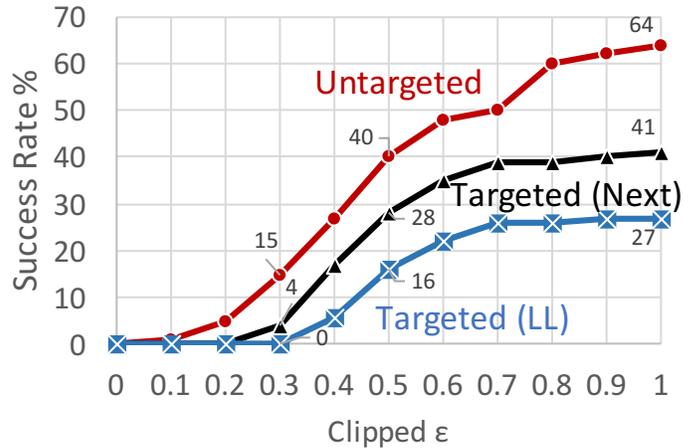


Fig. 7: Success rate on the adaptive adversarial examples. The success rate of the adaptive adversary decreases as we decrease the allowed ϵ value.

Nevertheless, the adaptive attack shows how an adversary may be able to find inputs that both confuse the model and are not detected as adversarial by the distance metrics. One obvious strategy for further complicating the adversary's task is to introduce randomness in the squeezing method. This is very different from attempts to obfuscate models, which have been shown vulnerable to transfer attacks. Instead, we can easily use cryptographic randomness to make the deployed framework unpredictable in ways that benefit the defender, since the adversary's search requires the knowledge on the exact squeezing operation. The defender has many opportunities to use randomness in selecting squeezing parameters (for example, instead of using a fixed 0.5 threshold for the 1-bit filter, using $0.5 \pm rand(0.1)$, or selecting random regions for the median smoothing instead of a fixed 2×2 region).

VI. CONCLUSION

The effectiveness of feature squeezing seems surprising since it is so simple and inexpensive compared to other proposed defenses. Developing a theory of adversarial examples remains an illusive goal, but our intuition is that the effectiveness of squeezing stems from how it reduces the search space of possible perturbations available to an adversary.

The bit depth reduction squeezers work essentially eliminate some of the lower bits, shrinking the feature space and forcing the adversary to produce larger perturbations. Since the features we effectively eliminate are not relevant for classification, this has little impact on the accuracy of legitimate samples. Bit depth reduction is effective in mitigating the adversarial examples generated by L_∞ attacks, and appears to be more effective against more advanced attacks since they use smaller perturbations. The spatial smoothing squeezers make pixels less different across an image, mitigating L_0 perturbations.

As discussed in Section V-D, feature squeezing is not immune to adversarial adaptation, but it substantially changes the challenge an adversary faces. Our general detection framework opens a new research direction in defending against adversarial examples and understanding the limits of deep neural networks in adversarial contexts.

⁵This work was done following public reports on the work in this paper; we shared details of our approach and code with the authors of [14], and much appreciate their sharing their implementation with us to enable the experiments reported here.

⁶He et al. reported results for CIFAR-10 with only bit depth reduction and median smoothing [14]. These results were similar to the results they reported on MNIST.

⁷In contrast, their target detector uses 0.3076 as threshold and uses a slightly different max function in combining dual-squeezers.

REFERENCES

- [1] Rodrigo Benenson. Classification datasets results. http://rodrigob.github.io/are_we_there_yet/build/classification_datasets_results.html.
- [2] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [3] Nicholas Carlini. Robust evasion attacks against neural network to find adversarial examples. https://github.com/carlini/nn_robust_attacks/.
- [4] Nicholas Carlini, Pratyush Mishra, Tavish Vaidya, Yuankai Zhang, Micah Sherr, Clay Shields, David Wagner, and Wencho Zhou. Hidden voice commands. In *USENIX Security Symposium*, 2016.
- [5] Nicholas Carlini and David Wagner. Defensive Distillation is not robust to adversarial examples. *arXiv preprint arXiv:1607.04311*, 2016.
- [6] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (Oakland)*, 2017.
- [7] Francois Chollet. Keras Implementation of Inception v3. https://github.com/fchollet/deep-learning-models/blob/master/inception_v3.py.
- [8] George E Dahl, Jack W Stokes, Li Deng, and Dong Yu. Large-scale malware classification using random projections and neural networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013.
- [9] Reuben Feinman, Ryan R Curtin, Saurabh Shintre, and Andrew B Gardner. Detecting adversarial samples from artifacts. *arXiv preprint arXiv:1703.00410*, 2017.
- [10] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2015.
- [11] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick McDaniel. On the (statistical) detection of adversarial examples. *arXiv preprint arXiv:1702.06280*, 2017.
- [12] Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples. *arXiv preprint arXiv:1412.5068*, 2014.
- [13] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and others. DeepSpeech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014.
- [14] Warren He, James Wei, Xinyun Chen, Nicholas Carlini, and Dawn Song. Adversarial example defenses: Ensembles of weak defenses are not strong. In *11th USENIX Workshop on Offensive Technologies (WOOT 17)*, 2017.
- [15] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- [16] Hossein Hosseini, Sreeram Kannan, Baosen Zhang, and Radha Poovendran. Deceiving google’s perspective api built for detecting toxic comments. *arXiv preprint arXiv:1702.08138*, 2017.
- [17] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [18] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten. Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993*, 2016.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- [20] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *International Conference on Learning Representations (ICLR) Workshop*, 2017.
- [21] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [22] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [23] Somshubra Majumdar. DenseNet Implementation in Keras. <https://github.com/titu1994/DenseNet/>.
- [24] Somshubra Majumdar. Keras Implementation of Mobile Networks. <https://github.com/titu1994/MobileNetworks/>.
- [25] Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. *arXiv preprint arXiv:1702.04267*, 2017.
- [26] Microsoft Corporation. Microsoft Malware Competition Challenge. <https://www.kaggle.com/c/malware-classification>, 2015.
- [27] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. <https://github.com/LTS4/universal/>.
- [28] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. DeepFool: a simple and accurate method to fool deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [29] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [30] Nicolas Papernot, Ian Goodfellow, Ryan Sheatsley, Reuben Feinman, and Patrick McDaniel. cleverhans v1.0.0: an adversarial machine learning library. *arXiv preprint arXiv:1610.00768*, 2016.
- [31] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against deep learning systems using adversarial examples. *arXiv preprint arXiv:1602.02697*, 2016.
- [32] Nicolas Papernot, Patrick 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.
- [33] Nicolas Papernot, Patrick McDaniel, Arunesh Sinha, and Michael Wellman. Towards the Science of Security and Privacy in Machine Learning. *arXiv preprint arXiv:1611.03814*, 2016.
- [34] Nicolas Papernot, Patrick 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 (Oakland)*, 2016.
- [35] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep Face Recognition. In *British Machine Vision Conference*, 2015.
- [36] Scientific Computing Tools for Python Developers. Multidimensional Image Processing ([scipy.ndimage](https://docs.scipy.org/doc/scipy/reference/tutorial/ndimage.html)). <https://docs.scipy.org/doc/scipy/reference/tutorial/ndimage.html>, 2009.
- [37] Scientific Computing Tools for Python Developers. Median Filter ([scipy.ndimage.median_filter](https://docs.scipy.org/doc/scipy/reference/generated/scipy.ndimage.median_filter.html#scipy.ndimage.median_filter)). https://docs.scipy.org/doc/scipy/reference/generated/scipy.ndimage.median_filter.html#scipy.ndimage.median_filter, 2017.
- [38] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- [39] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations (ICLR)*, 2014.
- [40] Matthew A Turk and Alex P Pentland. Face Recognition using Eigenfaces. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1991.
- [41] Beilun Wang, Ji Gao, and Yanjun Qi. A Theoretical Framework for Robustness of (Deep) Classifiers Under Adversarial Noise. In *International Conference on Learning Representations (ICLR) Workshop*, 2017.
- [42] Wikipedia. Median Filter. Page Version ID: 760708062.
- [43] Fei Zhang, Patrick PK Chan, Battista Biggio, Daniel S. Yeung, and Fabio Roli. Adversarial Feature Selection against Evasion Attacks. *IEEE Transactions on Cybernetics*, 2015.

TABLE VI: The non-default parameters used in this paper in generating adversarial examples. We didn't change the parameters for JSMA.

	Attack	Parameters	MNIST	CIFAR-10	ImageNet
L_∞	FGSM	eps	0.3	0.0156	0.0078
		BIM	eps	0.3	0.0080
		eps_iter	0.06	0.0012	0.002
	CW_∞	confidence	5		
L_2	DeepFool	overshoot	-	10	35
		confidence	5		
	CW_2	max_iterations	1000		
		batch_size	100	10	
L_0	CW_0	confidence	5		

APPENDIX

A. Attack Parameters

B. EvadeML-Zoo

The readers could reproduce the results using EvadeML-Zoo⁸.

⁸<https://github.com/mzweilin/EvadeML-Zoo>