# Divide, Denoise, and Defend against Adversarial Attacks

Seyed-Mohsen Moosavi-Dezfooli[1,2]

seyedmoosavi@apple.com

Ashish Shrivastava[1]

ashish.s@apple.com

Oncel Tuzel[1]

otuzel@apple.com

[1]Apple Inc., USA.    [2]École Polytechnique Fédérale de Lausanne, Switzerland.

## Abstract

*Deep neural networks, although shown to be a successful class of machine learning algorithms, are known to be extremely unstable to adversarial perturbations. Improving the robustness of neural networks against these attacks is important, especially for security-critical applications. To defend against such attacks, we propose dividing the input image into multiple patches, denoising each patch independently, and reconstructing the image, without losing significant image content. We call our method D3. This proposed defense mechanism is non-differentiable which makes it non-trivial for an adversary to apply gradient-based attacks. Moreover, we do not fine-tune the network with adversarial examples, making it more robust against unknown attacks. We present an analysis of the tradeoff between accuracy and robustness against adversarial attacks. We evaluate our method under black-box, grey-box, and white-box settings. On the ImageNet dataset, our method outperforms the state-of-the-art by 19.7% under grey-box setting, and performs comparably under black-box setting. For the white-box setting, the proposed method achieves 34.4% accuracy compared to the 0% reported in the recent works.*

## 1. Introduction

Deep neural networks (DNNs) have produced valuable results on many practical applications [28, 12, 51, 2, 26], but are vulnerable to even small adversarial perturbations [18, 39, 7]. In particular, such perturbations can change the decision of DNN-based image classifiers. The vulnerability of deep networks to adversarial manipulations of their input goes beyond classification tasks and additive perturbations [11, 34, 10, 25]. Moreover, the attacks are transferrable, meaning that an adversary can find these perturbations without having access to the network. For example, [29, 1] successfully attacked image classifiers used in commercial applications. These observations highlight the need to im-
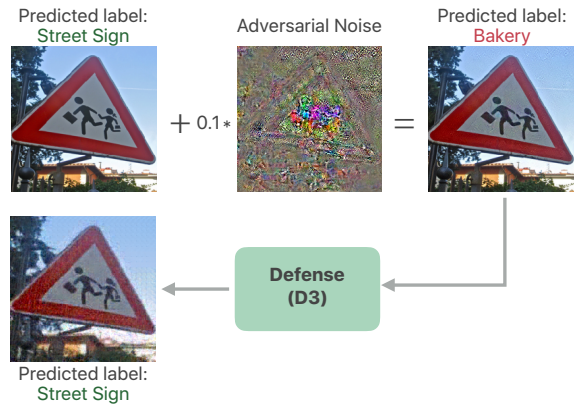


Figure 1: The proposed method (D3) transforms the input image using a non-differentiable algorithm. This transformation removes adversarial noise to improve robustness to attacks.

prove the robustness of deep networks, especially, if they are deployed in a hostile or security-critical environment.

Many existing defense methods either show results on small datasets, such as MNIST or CIFAR-10 [30, 6, 42], or they add adversarial examples to the training data [47]. Since new attacks are constantly being proposed, the defense should be attack-agnostic to make it robust against an unknown attack. An ideal defense should be non-differentiable so it does not allow the adversary to back-propagate through the defense mechanism.

Our defense mechanism (Figure 1) maps an input image to a new space and is based on the following observations: (1) increased dimensionality has an adverse effect on the robustness of deep networks [15], (2) the mapping should not reduce the accuracy of clean data, and (3) the mapping should be stable such that the output is minimally sensitive to input perturbations.

Our method divides the input image into multiple overlapping patches that are projected to a lower dimensional space using a dictionary. The dictionary consists of clean

image patches to reconstruct the input image patches with a variant of the matching pursuit (MP) algorithm [32]. We use this algorithm for denoising because it is fast and non-differentiable. We propose a novel patch-selection algorithm to construct the dictionary such that the selected patches are not too similar and they represent the salient parts of the image. This selection process mitigates the effect of adversarial perturbations while maintaining good accuracy on clean images.

We evaluate our algorithm using the ImageNet dataset under three settings – (1) black-box attacks where the adversary does not know about the network or the defense method, (2) grey-box attacks where the adversary knows the network parameters but not about the defense mechanism, and (3) white-box attacks where the adversary knows both the network parameters and the defense algorithm. Our algorithm performs comparably to other state-of-the-art methods on the black-box setting and performs significantly better on grey-box and white-box attacks. In addition, we show that as task complexity decreases (*e.g.* with a subset of the ImageNet classes), we can remove more information content from the image by denoising while maintaining high accuracy and robustness. Our contributions are:

- We propose a novel framework for defending against adversarial attacks by dividing images into overlapping patches and denoising them independently using a non-differentiable, attack-agnostic algorithm. Our patch-based method is particularly designed for high-dimensional datasets, e.g. ImageNet.

- We provide a thorough analysis of the tradeoff between clean image accuracy and robustness against adversarial attacks.

- We extensively evaluate our method on the ImageNet and CIFAR-10 datasets under different adversarial settings. We compare with the state-of-the-art papers and the NIPS-2017 challenge, and show improved performance.

## 2. Related Work

The seminal work by [46] highlights the vulnerability of deep neural networks to adversarial examples. Since then, many methods have been proposed to assess such vulnerability by developing various adversarial attacks. In [18], the authors proposed Fast Gradient Sign Method (FGSM) which attacks a classifier by computing the sign of the gradient of the loss w.r.t. the input images. To assess the robustness of deep networks more accurately, iterative algorithms such as DeepFool [37] and C&W [9] have later been introduced. It is also possible to use generative models like Generative Adversarial Networks (GANs) [17] to generate adversarial perturbations [5, 21].

In [46], it has also been shown that adversarial attacks are transferrable and can be used in black-box settings. In these settings, the adversary has access neither to the weights of the network nor to the architecture. More recently, it has been shown that there exists image-agnostic attacks, universal adversarial perturbation [36], which can be added to any image to fool a given network. Even worse, these perturbations can be computed without the dataset used for training the network [38].

In recent years, there have been several efforts to defend against adversarial attacks. The first defense against adversarial perturbations was proposed by [19] where they use stacked denoising auto-encoders to mitigate perturbations. A similar approach has been studied in [33] to denoise adversarial examples. Recently, other generative models, e.g. GANs and PixelCNN [48], have been used to project back the malicious samples on the manifold of data [42, 45]. However, such methods are restricted to small-scale datasets such as MNIST and CIFAR-10. In [40], distillation was suggested as a defense; however, it only masks the gradient and is still vulnerable in black-box settings as demonstrated in [9]. In [30], authors provide an efficient adversarial training framework based on the robust optimisation to counter the first-order adversaries. However, their method is not model-agnostic and due to computational complexities, it has not been applied on large-scale datasets such as ImageNet. Recently, in [47], adversarial training has been applied on ImageNet using an ensemble of networks. The main drawback of such adversarial training scheme is that it overfits the noise and does not generalize well against an unknown attack [20]. Very recently, several strategies to defend against attacks have been explored using various image transformations [20]. As a different approach, detecting malicious samples, instead of improving the robustness, is sought in [35, 16]. They demonstrated that deep networks can be augmented with a network to detect adversarial examples. This approach also suffers from overfitting to specific types of perturbations. Recent work, [8], has successfully attacked ten different defense strategies emphasizing the difficulty of this problem.

There have been few theoretical works studying the robustness of deep networks. There is a tradeoff between accuracy and robustness of kernel classifiers [14]. In [15], the authors established a bound on the robustness of a certain type of classifiers when the adversary is restricted to a low dimensional space. In [44, 23], some lower bounds have been derived on the robustness of simple neural networks.

## 3. Problem Formulation

Let $f_{\boldsymbol{\theta}}(\boldsymbol{x}) : \mathbb{R}^d \to \mathbb{N}$ be a classifier parameterized by $\boldsymbol{\theta}$ that computes the class of an input image $\boldsymbol{x}$, where $\mathbb{N}$ is the set of natural numbers denoting class labels. An adversary can perturb the image with noise $\boldsymbol{v}$ such that $f_{\boldsymbol{\theta}}(\boldsymbol{x}) \neq f_{\boldsymbol{\theta}}(\boldsymbol{x} + \boldsymbol{v})$. The robustness of the classifier at
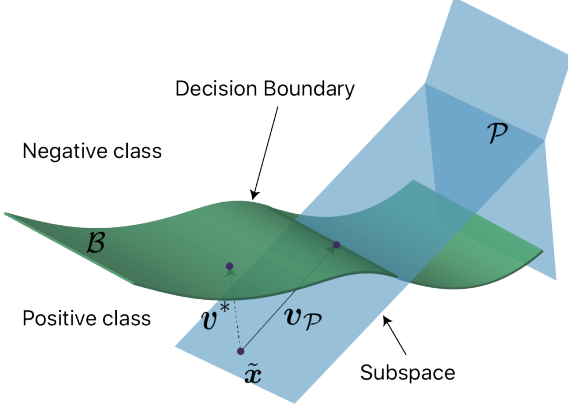
Figure 2: Reducing dimensionality improves the robustness of the classifier. $\mathcal{P}$ is a union of subspaces illustrated by the blue hyper-planes. Here, $\tilde{x} = T(x)$ is the image projected to the nearest subspace in $\mathcal{P}$. To avoid clutter, we are not showing the original image $x$. The adversarial noise $v^*$ is the smallest distance from $\tilde{x}$ to the classifier's decision surface $\mathcal{B}$. When the adversary is restricted to a smaller dimensional subspace (the projected hyper-plane), the norm of the noise $v_{\mathcal{P}}$ is much bigger than the norm of $v^*$ to cross the decision boundary.
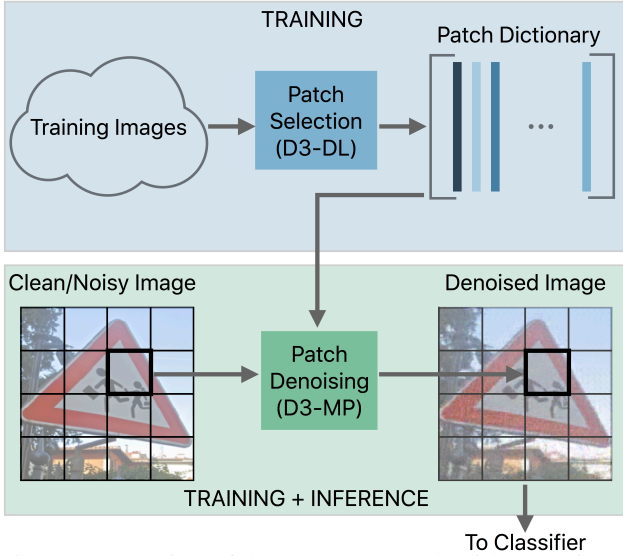


Figure 3: Overview of the proposed D3 algorithm. An input image is divided into patches and each patch is denoised independently.

$x_0$, denoted by $\rho(x_0)$, can be defined as the minimum perturbation needed to change the predicted label [15]:

$$\rho(x_0) = \arg\min_{v} \|v\|_2, \text{ s. t. } f_\theta(x) \neq f_\theta(x + v). \quad (1)$$

The noise can be scaled to make the attack stronger. We improve robustnesss by learning a stable transformation, $T(.)$, such that the label of the transformed image does not

change when corrupted with the noise, *i.e.* $f_\theta(T(x)) = f_\theta(T(x + v))$. Moreover, we want to maintain the accuracy when the clean images are transformed by $T(.)$ ensuring that the $f_\theta(T(x))$ is equal to the ground truth label of $x$. Most attacks rely on computing gradients of the classification function w.r.t. the input. Hence, it is desirable to make the transformation $T(.)$ non-differentiable so that the gradients cannot pass through the defense block. To create a defense mechanism that is robust against an unknown future attack, we want to keep the defense algorithm to be attack-agnostic and do not want to fine-tune the network with simulated adversarial images.

Under some regularity conditions, restricting the adversary to a low dimensional subspace can improve robustness [15]. Simple dimensionality reduction methods, such as PCA, have been studied in [6, 24] to improve defense against adversarial perturbations. However, such solutions only work on simpler tasks (such as classification on MNIST digits) as they significantly decrease the discriminative performance of the network. Such transformations usually remove the high frequency information required for complex tasks such as 1000-class classification on ImageNet dataset [12]. Therefore, a better information preserving dimensionality reduction method is required to limit the space of adversarial noise, while keeping important details.

## 4. The D3 Algorithm

Assume that an operator $T(x)$ projects the input image $x \in \mathbb{R}^d$ to the closest subspace in a union of $m$ dimensional subspaces. This operation is a linear projection operator onto an $m$ dimensional subspace in a local neighborhood of $x$. For additive perturbations, the adversary is limited to *locally* seeking noise in an $m$ dimensional subspace and the robustness, $\rho$, can ideally be improved by a factor of $\sqrt{d/m}$ [15]. This intuition is illustrated in Figure 2. Motivated by this result, we look for a transformation, $T(.)$, satisfying the following conditions: (1) rank $J_T(x) \ll d$, where $x \in \mathbb{R}^d$, where $J_T(x)$ is the Jacobian matrix in a small neighborhood of $x$, and (2) $\|T(x) - x\|$ is small. The first condition ensures that the local dimensionality is low, and the second condition means that the image and its transformed version are close enough to each other to preserve visual information.

We propose a patch-based denoising method for defense. We divide the input image into multiple patches, and denoise them independently with sparse reconstruction using a dictionary of patches. Assume that $\kappa$ is the sparsity (number of components used to reconstruct a patch) and each patch is $P \times P$ pixels. Then, for non-overlapping patches, the local dimensionality of our projection operator $T(.)$ would be $\kappa \frac{d}{P^2}$. According to [15], this dimensionality reduction would ideally improve robustness by a factor of $\frac{P}{\sqrt{\kappa}}$.

Sparse reconstruction and dictionary-based methods have

been widely used to enhance the quality of images [3, 31, 13, 49, 43, 41]. For computational efficiency, we use sampled image patches as our dictionary. We use a novel patch selection algorithm that is optimized to improve robustness of the classifier. For sparse reconstruction, we use an efficient greedy algorithm which is a variant of matching pursuit [32] . Our method is summarized in Figure 3.

Computing the performance of the D3 algorithm for different hyper-parameters requires both training the D3 algorithm and the classification network, which are computationally expensive. Therefore, to efficiently study the effect of hyper-parameters of our algorithm, we compute the following proxy metrics:
(1) **"Matching-rate"** (MR) is the fraction of patches that are identical in the denoised image $T(x + v)$ and the clean image $T(x)$. Let $\{p_1, p_2, \ldots, p_n\}$ and $\{\hat{p}_1, \hat{p}_2, \ldots, \hat{p}_n\}$ be patches extracted from $x$ and $x + v$, respectively. The matching-rate is defined as, MR $= \mathbb{E}_{x \in \mathcal{D}}(\gamma(x))$, where

$$\gamma(x) = \frac{1}{n} \sum_{j=1}^{n} \mathbb{1}_{\|T(p_j) - T(\hat{p}_j)\|_{\infty} \leq \delta}$$

and $\mathbb{1}_{[.]}$ is an indicator function. Here, with slight abuse of notation, we assume that $T(.)$ is applied to patches. Higher MR corresponds to being more robust to the attacks.
(2)**"Reconstruction-error"** (RE) is the average $\ell_2$ distance between the clean image, $x$, and the transformed image, $T(x)$: RE $= \mathbb{E}_{x \in \mathcal{D}}(\|x - T(x)\|_2 / \|x\|_2)$. A higher reconstruction-quality, $(1 - \text{RE})$, results in higher classification accuracy for the clean images as more information is retained.

Our experiments show that these proxy metrics are highly correlated with accuracy and robustness of the classifier. We use 500 images randomly chosen from the dataset to quickly compute these values. Next, we describe the patch denoising algorithm (D3-MP), and the patch-selection algorithm to construct the dictionaries (D3-DL).

### 4.1. Patch-based Denoising (D3-MP)

Let $\{S_i\}_{i=1}^{\kappa}$ be a set of dictionaries computed using our patch-selection algorithm. Each dictionary $S_i \in \mathbb{R}^{P^2 \times \eta}$ is a matrix containing $\eta$ columns of dimension $P^2$. Here, $\kappa$ is the sparsity. The first dictionary $S_1$ is used to select the first atom, while reconstructing a given patch $p$. Then, the residual is computed between the image patch, $p$, and the selected atom, $s_l$. As in the standard matching pursuit (MP), the residual is used to select the next atom. But unlike standard MP, we use a different dictionary, $S_i$, to select at the $i^{\text{th}}$ sparsity level. We provide pseudo-code for this approach in Algorithm 1.

### 4.2. Patch Selection Algorithm to Learn Dictionaries (D3-DL)

To scale up the dictionary learning task for a large dataset such as the ImageNet, we propose an efficient greedy patch-selection algorithm. As mentioned earlier, we compute multiple dictionaries for different sparsity levels in our D3-MP algorithm.

We build the set of dictionaries in a greedy manner by selecting the "important" and "diverse" patches. The algorithm takes into account the saliency information of images. The norm of the gradient of the classification function w.r.t. to the input image is used as the saliency map. We do importance sampling among all the patches w.r.t. the saliency map. We add a patch to the dictionary if the reconstruction of the patch using the existing dictionary has greater than a threshold angular distance, $\epsilon$, from the patch. The saliency map helps to preserve the details that are important for the classification task, and the cutoff on the angular distance ensures that the dictionary is diverse.

In our experiments, using a pre-tained network on the ImageNet dataset, we find $4\%$ improvement in classification accuracy with the saliency map compared to randomly selecting a patch from the whole image. The diversity among dictionary atoms encourages mapping a clean and corresponding noisy image patch to the same dictionary atom. Ensuring that any two patches from the dictionary are a certain threshold apart also improves the MR and the robustness of the classifier.

After the first dictionary is constructed, we reconstruct the image patches using this dictionary and compute the residuals. The next dictionary is constructed on the residual images instead of the original images. This process is repeated for all the remaining dictionaries (see Algorithm 2 for pseudo-code). We found that the MR and 1 - RE were higher when we used a different dictionary for different sparsity levels using residual images compared to using one common dictionary for all sparsity levels constructed from original images. For example, with $\kappa = 2$, we found the MR = 0.88, 1 - RE = 0.83 when using two separate dictionaries where the second dictionary contained residuals instead of image patches. Comparing those to MR = 0.80, 1 - RE = 0.81 when using one dictionary, constructed using original image patches, shows that using multiple dictionaries gives better matching-rate and reconstruction quality.

### 4.3. Denoising Algorithm

The proposed defense algorithm (D3): (1) divides the input image into overlapping patches, (2) denoises each patch (with D3-MP) using the constructed dictionaries (with D3-DL), and (3) reconstructs the denoised image by averaging the pixels in overlapping patches. In our experiments, we set the amount of overlap to $75\%$ of the patch size.
**Randomization:** In our experiments, we observe that

**Algorithm 1** D3-MP

**Input:** A set of dictionaries $\{S_i\}_{i=1}^{\kappa}$, image patch $p$.
**Output:** processed patch $q$.
$q \leftarrow 0$
$\hat{p} \leftarrow p$
**for** $i = 1$ **to** $\kappa$ **do**
    $a \leftarrow \hat{p}^{\top} S_i$
    $l \leftarrow \mathbf{argmax}_k |a_k|$
    $q \leftarrow q + a_l s_l$
    $\hat{p} \leftarrow \hat{p} - a_l s_l$
**return** $q$

---

**Algorithm 2** D3-DL

**Input:** saliency algorithm $\mathcal{H}$, training images $\mathcal{D}$, size of dictionary $\eta$, sparsity $\kappa$, $\epsilon$.
**Output:** set of dictionaries $\{S_i\}_{i=1}^{\kappa}$.
**for** $i = 1$ **to** $\kappa$ **do**
    $n \leftarrow 0$
    $S_i \leftarrow []$
    **while** $n < \eta$ **do**
        Randomly select $x \in \mathcal{D}$.
        Compute saliency map $\mathcal{H}(x)$.
        Randomly select patch $s$ from $x$ according to $\mathcal{H}(x)$.
        **if** $i = 1$ **then**
            $\tilde{s} \leftarrow$ D3-MP$(\{S_1\}, s)$
            //Add $s$ to $S_1$ if it is arcsin($\epsilon$) away.
            **if** $\|s - \tilde{s}\|_2 / \|s\|_2 > \epsilon$ **then**
                Concatenate $s/\|s\|_2$ to columns of $S_i$
                $n \leftarrow n + 1$
        **else**
            $\tilde{s} \leftarrow$ D3-MP$(\{S_j\}_{j=1}^{i-1}, s)$
            $r \leftarrow s - \tilde{s}$
            //Add $r$ to $S_i$ if it is arcsin($\epsilon$) away.
            $\tilde{r} \leftarrow$ D3-MP$(\{S_i\}, r)$
            **if** $\|r - \tilde{r}\|_2 / \|r\|_2 > \epsilon$ **then**
                Concatenate $r/\|r\|_2$ to columns of $S_i$
                $n \leftarrow n + 1$

---

without giving access to the patch dictionaries (*i.e.* the grey-box setting), the D3 algorithm is successful in defending against the adversarial attack. However, the defense is weaker when the adversary has access to the dictionary. This observation motivates us to add randomization to our transformation function, $T(.)$, when the adversary has full access to the D3 algorithm (white-box setting). By adding randomization, even though the adversary can access the dictionary, exact atoms used for reconstruction will not be available. We add following efficient randomization schemes (both training and test time) to our defense:

- We randomize over the columns of the dictionaries by randomly selecting one fifth of the atoms in the patch

dictionaries while denoising a patch.

- We further randomize by first selecting the top-2 most correlated atoms from the patch dictionary and randomly picking one of those two.
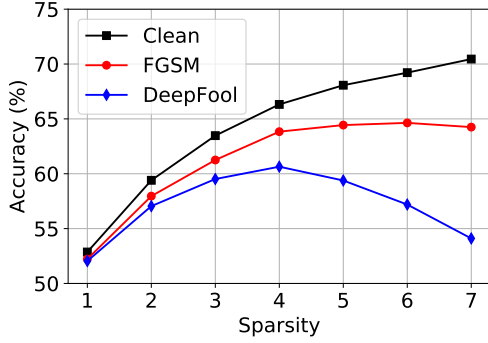
## 5. Experiments

We conduct experiments on the 1000-class ImageNet dataset [12] with a ResNet-50 [22]. We evaluate defense against the following diverse set of attack algorithms – **FGSM** [18] which is a one-step attack where gradient direction is multiplied by the norm of the noise, **DeepFool** [37] that is an iterative attack ensuring that the network output is changed, **CWattack** [9] that solves an optimization problem to find the adversarial perturbation, and **UAP** [36] as a transferable attack that computes a universal noise for all the images. We compute the adversarial perturbation in floating point and scale it to have a fixed $\ell_2$ norm, *i.e.* $\frac{\|v\|_2}{\|x\|_2} = 0.06$. This noise is added to the image after it has been converted to floating point and normalized. This experimental setup is described in [20]. We compare the D3 algorithm with several image transformation based defenses proposed in [20] and show significant improvement in classification accuracy. Furthermore, we analyze the tradeoff between the clean image accuracy and the robustness against the attacks. We study the effect of the hyper-parameter values and show how they can be tuned to improve the accuracy or the robustness. We use a pre-trained network on the original image space $x$ and fine-tune it on the transformed images $T(x)$. We set $\epsilon = 0.85$ for dictionary reconstruction and discuss the effects of other hyper-parameters in Section 5.1.
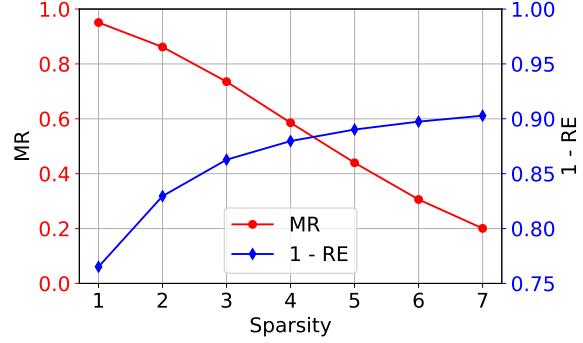
**Tradeoff between Accuracy and Robustness:** Hyper-parameters of the D3 algorithm can be used to control the quality of image reconstruction. For example, as we increase the sparsity, $\kappa$, more details are preserved and the accuracy on the reconstructed clean images, $T(x)$, is improved. However, if the image has been corrupted with adversarial noise, increasing $\kappa$ reconstructs more noise, and the accuracy on the reconstructed noisy images, $T(x + v)$, eventually decreases.

We plot the accuracy on clean and noisy images with a grey-box attack in Figure 4(a). As we can see from the plots, the classification accuracy on the clean images consistently improves with sparsity. The accuracy on the noisy images is low with small sparsity and increases in the beginning as the reconstruction quality improves. After increasing the sparsity to a large value, the noise also starts getting reconstructed, and the classification accuracy drops. Hence, there is a "sweet spot" (for example, $\kappa = 4$ for DeepFool) for optimal defense against the attack, and sparsity can be used as a tradeoff parameter.

Figure 4(b) shows the matching-rate (MR) and the reconstruction-quality (1 - RE) for the same sparsity val-

Figure 4: (a) Tradeoff between the clean image accuracy and the robustness of the classifier with increasing sparsity. (b) Matchine rate (MR) and reconstruction quality (1 - RE) for the same sparsity values. From these two figures, we can see that clean image accuracy and the reconstruction quality increase with sparsity. The MR decreases with sparsity, causing the classifier to be less robust. The classifier performance on adversarially perturbed images is low with small sparsity due to low reconstruction quality, and is also low with large sparsity because of low matching-rate. The optimal performance is achieved at an intermediate value (*e.g.* $\kappa = 4$ for DeepFool attack). Due to computational constraints, we reduce the patch overlap to 8 pixels for this analysis which improves the reconstruction speed by 4x.

ues. We see that the reconstruction quality improves with sparsity, and is correlated with the classification accuracy on the clean images. However, the matching-rate decreases as sparsity is increased, causing the network to be less robust to adversarial noise. In our experiments, we find that dictionary size also plays a role in the accuracy and robustness tradeoff. A larger dictionary improves the accuracy on the clean images because the images are better reconstructed. A smaller dictionary generally improves the robustness because the dictionary atoms are, on average, farther apart. This correlation between the (MR, 1-RE) and the classifier's performance enables us to efficiently study the effects of hyper-parameters.

Based on this analysis, we pick three settings to evaluate our algorithm. In the first setting, we set high sparsity and choose a large dictionary ($\kappa = 5$, dictionary size $\eta = 40k$). This setting encourages high accuracy on the clean images, but is less robust against adversarial attacks. In the second setting, we reduce the dictionary size while keeping the sparsity the same ($\kappa = 5$, dictionary size $\eta = 10k$). In the third setting, we set the dictionary size $\eta = 10k$, and reduce the sparsity to $4$, further improving the robustness. Next, we describe three types of attacks (black-box, grey-box, and white-box) and compare our results with the state-of-the-art.

**Black-box Attack:** In this attack, the adversary does not have access to the network parameters, or the defense mechanism. We use a separate ResNet-50 network and compute the adversarial perturbations using original images, $x$. As mentioned earlier, the $\ell_2$-norm of the noise is set to 0.06 and it is added to the floating point image. As shown in

Table 1, all the methods are robust to these attacks as this attack is the easiest to defend. Note that here we assume that the adversary knows about the network architecture but not the exact parameters. We also compare with the best method [27] in the *NIPS2017 challenge*. This method used a different base network (InceptionV3) which had a higher clean image (no-attack) accuracy of $76.6\%$. We used their pre-trained network/defense algorithm and evaluated it with our experimental setup. Compared to D3, it had a significantly lower accuracy under DeepFool attack ($59.8\%$) and slightly better accuracy ($71.9\%$) under FGSM attack. Note that better performance on FGSM attack may be due to their defense being trained using adversarially perturbed images with FGSM.

**Grey-box Attack:** In this setting, the adversary does not have access to the defense mechanism but knows about the network weights. We compute all the noise patterns using gradients of the fine-tuned network evaluated at the original images, $x$. This setting is the same as the "white-box" setting in [20], so we compare our results in grey-box setting to their white-box setting. Without any defense mechanism, the attacks are very successful in this setting: FGSM reduces the classification accuracy to $6.2\%$, DeepFool reduces it to $9.2\%$, and UAP reduces it to $20.8\%$. Our results in Table 2 show that we perform significantly better than the state-of-the-art. We also note that as we decease the dictionary size $\eta$ or the sparsity $\kappa$, the robustness of the classifier improves. For example, with $\eta = 40k, \kappa = 5$, the DeepFool accuracy is $58.7\%$ which improves to $64.4\%$ by reducing $\eta$ to $10k$ and sparsity to $4$. As before, we do not use any adversarial

Table 1: Top-1 classification accuracy (%) with black-box attacks (1000 ImageNet classes).

| METHOD | NO ATTACK | DEEPFOOL | CWATTACK | FGSM | UAP |
|---|---|---|---|---|---|
| D3 ($\eta = 40k, \kappa = 5$) | 71.8 | 63.1 | 63.3 | **68.6** | **71.5** |
| D3 ($\eta = 10k, \kappa = 5$) | 70.8 | 64.6 | 64.8 | 68.3 | 70.3 |
| D3 ($\eta = 10k, \kappa = 4$) | 69.0 | 64.8 | 65.0 | 67.1 | 68.9 |
| QUILTING [20] | 70.1 | 65.2 | 64.1 | 65.5 | - |
| TVM + QUILTING [20] | **72.4** | 65.8 | 64.0 | 65.7 | - |
| CROPPING + TVM + QUILTING [20] | 72.1 | **67.1** | **65.3** | 66.7 | - |

Table 2: Top-1 classification accuracy (%) with grey-box attacks (1000 ImageNet classes).

| NETWORKS | NO ATTACK | DEEPFOOL | CWATTACK | FGSM | UAP |
|---|---|---|---|---|---|
| D3 ($\eta = 40k, \kappa = 5$) | 71.8 | 58.7 | 57.9 | 67.3 | **71.5** |
| D3 ($\eta = 10k, \kappa = 5$) | 70.8 | 62.3 | 62.4 | 67.5 | 70.4 |
| D3 ($\eta = 10k, \kappa = 4$) | 69.0 | **64.4** | **64.5** | 67.1 | 68.7 |
| QUILTING [20] | 66.9 | 34.5 | 30.5 | 39.6 | - |
| CROPPING [20] | 65.4 | 44.9 | 41.1 | 49.5 | - |
| TVM [20] | 66.3 | 44.7 | 48.4 | 31.4 | - |
| ENSEMBLE TRAINING [47, 20][1] | **80.3** | 1.8 | 22.2 | **69.2** | - |

Table 3: Top-1 classification accuracy (%) with white-box attacks (1000 ImageNet classes).

| NETWORKS | NO ATTACK | DEEPFOOL | CWATTACK | FGSM |
|---|---|---|---|---|
| XIE ET AL. [50, 4][2] | - | - | 0.0 | - |
| QUILTING [20, 4][2] | - | - | 0.0 | - |
| D3 ($\eta = 40k, \kappa = 5$) | **70.8** | 27.5 | 27.1 | 54.5 |
| D3 ($\eta = 10k, \kappa = 5$) | 69.8 | 31.0 | 30.9 | 55.8 |
| D3 ($\eta = 10k, \kappa = 4$) | 68.2 | **34.3** | **34.4** | **56.9** |

examples while fine-tuning the network. Ensemble training [47][1] does better on FGSM attack (69.2% accuracy) by training the network with FGSM adversarial examples. But the trained model performs poorly with DeepFool attack, resulting only in 1.8% accuracy.

**White-box Attack:** Since the D3 transformation function is not differentiable, we cannot fool the network using gradient-based attacks (FGSM, DeepFool, CWAttack, or UAP). We compute the adversarial noise, $v$, using the fine-tuned network weights while the gradients are computed at the transformed image, $T(x)$. This noise, $v$, is added to the image $x$. This attack is same as BDPA of [4]. Under this challenging setting, the DeepFool attack reduces the classification accuracy to 13.0%, and the FGSM attack reduces it to 34.4%. We can make D3 more robust by adding randomization, which prevents the attacker from accessing the exact atoms used for reconstruction, as described in Section 4.3. Classification accuracies using D3 with randomization are

shown in Table 3 which shows significant improvement over the deterministic version. As with the other attack types, the algorithm is more robust with a smaller dictionary and lower sparsity. We achieve accuracy of 13% without randomization and 34.4% with randomization on the BPDA compared to the 0% accuracy reported in [4]. We hypothesize that this improvement is partially due to the fact that our transformation radically reduces the effective dimensionality of the input images. We are able to use a bigger patch size $16 \times 16$ (even upto $32 \times 32$) which effectively limits the search space for an adversary.

## 5.1. Hyper-parameters selection

The proposed D3 defense has the following hyper-parameters: patch-size, $P$, sparsity, $\kappa$, dictionary size, $\eta$, and minimum distance between two dictionary atoms, $\epsilon$. They affect the overall performance in different ways. For example, as we increase the patch size, the robustness of the classifier improves while it becomes less accurate. In-

---

[1]The results of the ensemble training algorithm are taken from [20].
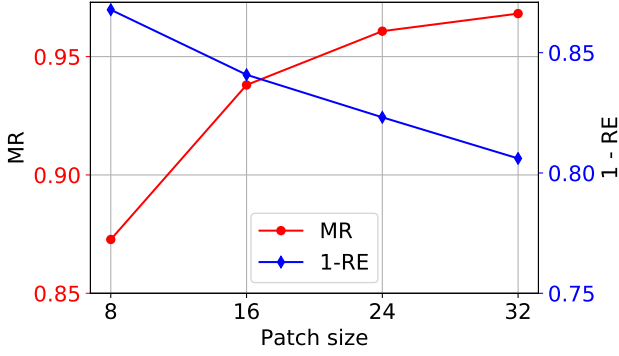
[2]The results of the BPDA attack are taken from [4].

Figure 5: Effect of patch size ($P$) on matching-rate and reconstruction-quality (1 - RE). Increasing $P$ improves the matching-rate (robustness) but reduces the reconstruction-quality.
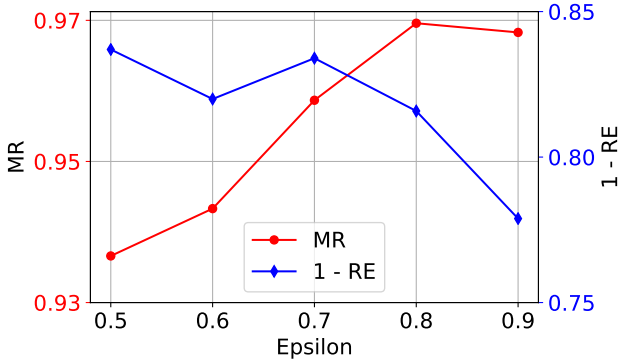


Figure 6: Matching-rate improves as $\epsilon$ increases or correlation among dictionary atoms decreases. The reconstruction-quality ($1 - \text{RE}$) decreases with increasing $\epsilon$.

creasing the minimum angle between two atoms, i.e. $\epsilon$, also improves the robustness but degrades the classifier accuracy. However, we recommend keeping $\epsilon$ high because it improves the robustness more than it hurts accuracy on the clean images.

We analyze the effect of $\kappa$ on accuracy in Figure 4, and show how the accuracy and the robustness are correlated with the matching-rate (MR) and the reconstruction-quality (1 - RE). We show the effect of different patch sizes ($P = 8, 16, 24, 32$) in Figure 5 on matching-rate and reconstruction-quality. As we increase the patch size, the matching-rate increases that is the robustness of the classifier improves. However, the reconstruction-quality ($1 - \text{RE}$) decreases making the classifier less accurate.

Next, we study the effect of $\epsilon$ in Figure 6. As expected, increasing the minimum angle between two atoms, increases the matching-rate and will improve the robustness. We also find that the reconstruction-quality ($1 - \text{RE}$) decreases, thus decreasing the classifier accuracy. Figure 7 shows example images with different sparsity values.
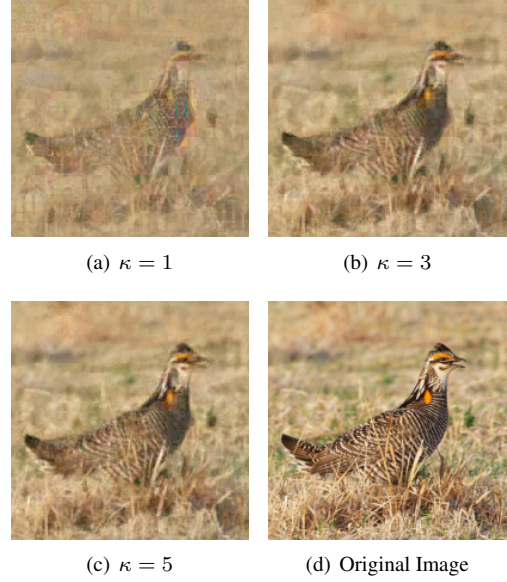


Figure 7: Denoised images with varying sparsity.

**Effect of Task-Complexity:** For simpler tasks, the D3 algorithm can reduce more information from the images while maintaining the classification accuracy and the robustness to attack. To evaluate our defense mechanism on a lower complexity task, we train a classifier on 50 randomly selected ImageNet classes. Under black-box attack, in the worst case the performance reduces from 93.2% to 91.8% and under grey-box attack it reduces to 89.7%. In the white-box setting, the top-1 classification accuracy on adversarial images drops to 70.9% from 91.7% (on clean images). We also evaluated our algorithm on CIFAR-10 dataset, another lower complexity task, and compared with many state-of-the-art methods given in [45] under FGSM attack. Our algorithm achieved 87% clean accuracy and 80% accuracy under the attack. Our attack-agnostic defense is better than all the defense methods but the one that specialized specifically for FGSM attack by adding adversarial images to the training. The complete evaluation on small ImageNet and CIFAR-10 datasets are included in Appendices A and B, respectively.

## 6. Conclusion

We described a novel patch-based denoising algorithm to improve the robustness of classifiers, trained on large-scale datasets, against adversarial perturbations. We developed two proxy metrics (MR and RE) to help guide us in designing the algorithm. The design of our denoising (D3-MP), and patch selection (D3-DL) algorithms encourages high matching-rate between the clean patch and the corresponding noisy patch. We provided a thorough study of the tradeoff between clean image accuracy and robustness against the attacks. We proposed an efficient randomization schemes to improve the robustness against white-box attacks. Our eval-

uation on the ImageNet show that our defense mechanism provides state-of-the-art results.

# References

[1] Wieland Brendel *, Jonas Rauber *, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *Proc. ICLR*, 2018. 1

[2] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016. 1

[3] M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, Nov 2006. 4

[4] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *Proc. ICML*, 2018. 7

[5] Shumeet Baluja and Ian Fischer. Adversarial transformation networks: Learning to generate adversarial examples. *arXiv preprint arXiv:1703.09387*, 2017. 2

[6] Arjun Nitin Bhagoji, Daniel Cullina, and Prateek Mittal. Enhancing robustness of machine learning systems via data transformations. *arXiv preprint arXiv:1704.02654*, 2017. 1, 3

[7] Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1704.02654*, 2017. 1

[8] Nicholas Carlini and David A. Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *ACM Workshop on Artificial Intelligence and Security*, 2017. 2

[9] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 2017. 2, 5

[10] Bo Li Warren He Mingyan Liu Dawn Song Chaowei Xiao, Jun-Yan Zhu. Spatially transformed adversarial examples. In *Proc. ICLR*, 2018. 1

[11] Moustapha Cissé, Yossi Adi, Natalia Neverova, and Joseph Keshet. Houdini: Fooling deep structured visual and speech recognition models with adversarial examples. In *Proc. NIPS*, 2017. 1

[12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *Proc. CVPR*, 2009. 1, 3, 5

[13] Michael Elad. *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing*. Springer, 1st edition, 2010. 4

[14] Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Analysis of classifiers' robustness to adversarial perturbations. *arXiv preprint arXiv:1502.02590*, 2015. 2

[15] Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Robustness of classifiers: from adversarial to random noise. In *Proc. NIPS*, 2016. 1, 2, 3

[16] Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner. Detecting adversarial samples from artifacts. *arXiv preprint arXiv:1703.00410*, 2017. 2

[17] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Proc. NIPS*, 2014. 2

[18] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proc. ICLR*, 2015. 1, 2, 5

[19] Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples. *arXiv preprint arXiv:1412.5068*, 2014. 2

[20] Chuan Guo, Mayank Rana, Moustapha Cissé, and Laurens van der Maaten. Countering adversarial images using input transformations. In *Proc. ICLR*, 2018. 2, 5, 6, 7

[21] Jamie Hayes and George Danezis. Machine learning as an adversarial service: Learning black-box adversarial examples. *arXiv preprint arXiv:1708.05207*, 2017. 2

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016. 5

[23] Matthias Hein and Maksym Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. *arXiv preprint arXiv:1705.08475*, 2017. 2

[24] Dan Hendrycks and Kevin Gimpel. Visible progress on adversarial images and a new saliency map. *arXiv preprint arXiv:1608.00530*, 2016. 3

[25] Can Kanbak, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Geometric robustness of deep networks: analysis and improvement. *arXiv preprint arXiv:1711.09115*, 2017. 1

[26] Ivan Krasin, Tom Duerig, Neil Alldrin, Andreas Veit, Sami Abu-El-Haija, Serge Belongie, David Cai, Zheyun Feng, Vittorio Ferrari, Victor Gomes, Abhinav Gupta, Dhyanesh Narayanan, Chen Sun, Gal Chechik, and Kevin Murphy. OpenImages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from https://github.com/openimages*, 2016. 1

[27] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Jun Zhu, and Xiaolin Hu. Defense against adversarial attacks using high-level representation guided denoiser. *arXiv preprint arXiv:1712.02976*, 2017. 6

[28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Proc. ECCV*, 2014. 1

[29] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *Proc. ICLR*, 2016. 1

[30] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *Proc. ICLR*, 2018. 1, 2

[31] J. Mairal, M. Elad, and G. Sapiro. Sparse representation for color image restoration. *IEEE Transactions on Image Processing*, Jan 2008. 4

[32] S.G. Mallat and Zhifeng Zhang. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on Signal Processing*, Dec 1993. 2, 4

[33] Dongyu Meng and Hao Chen. Magnet: A two-pronged defense against adversarial examples. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, CCS '17, pages 135–147, New York, NY, USA, 2017. ACM. 2

[34] J. H. Metzen, M. C. Kumar, T. Brox, and V. Fischer. Universal adversarial perturbations against semantic image segmentation. In *Proc. ICCV*, 2017. 1

[35] Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. In *Proc. ICLR*, 2017. 2

[36] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Proc. CVPR*, 2017. 2, 5

[37] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *Proc. CVPR*, 2016. 2, 5

[38] Konda Reddy Mopuri, Utsav Garg, and R Venkatesh Babu. Fast feature fool: A data independent approach to universal adversarial perturbations. In *Proc. BMVC*, 2017. 2

[39] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple black-box adversarial perturbations for deep networks. *arXiv preprint arXiv:1612.06299*, 2016. 1

[40] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy*, 2016. 2

[41] Y. C. Pati, R. Rezaiifar, and P. S. Krishnaprasad. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Asilomar Conference on Signals, Systems and Computers*, 1993. 4

[42] Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *Proc. ICLR*, 2018. 1, 2

[43] Ashish Shrivastava, Vishal M. Patel, Jaishanker K. Pillai, and Rama Chellappa. Generalized dictionaries for multiple instance learning. *International Journal of Computer Vision*, Sep 2015. 4

[44] Aman Sinha, Hongseok Namkoong, and John C. Duchi. Certifiable distributional robustness with principled adversarial training. In *Proc. ICLR*, 2018. 2

[45] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *Proc. ICLR*, 2018. 2, 8, 11

[46] Christian Szegedy, Google Inc, Wojciech Zaremba, Ilya Sutskever, Google Inc, Joan Bruna, Dumitru Erhan, Google Inc, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *Proc. ICLR*, 2014. 2

[47] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In *Proc. ICLR*, 2018. 1, 2, 7

[48] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *Proc. ICML*, 2016. 2

[49] J. Wright, Y. Ma, J. Mairal, G. Sapiro, T. S. Huang, and S. Yan. Sparse representation for computer vision and pattern recognition. *Proceedings of the IEEE*, June 2010. 4

[50] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan L. Yuille. Mitigating adversarial effects through randomization. In *Proc. ICLR*, 2018. 7

[51] Yuting Zhang, Kibok Lee, and Honglak Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. In *Proc. ICML*, 2016. 1

## Appendix A: Experiments on 50-class ImageNet

For simpler tasks, the D3 algorithm can reduce more information from the images while maintaining the classification accuracy and the robustness to attack. To evaluate our defense mechanism on a lower complexity task, we train a classifier on 50 randomly selected ImageNet classes. We find that the proposed defense (D3) is successful on this lower complexity task.

As shown in Tables 4 and 5, black and grey-box attacks lower the top-1 classification accuracy by less than 4%. In the most challenging white-box setting including the DeepFool attack, the top-1 classification accuracy drops to 70.9% from 91.7% on clean images. This 20.8% drop compares favorably with the 35.3% loss in accuracy for the 1000-class classifier under the same setting.

Table 4: Top-1 classification accuracy (%) on black-box attacks (50 ImageNet classes)

| NETWORKS | NO ATTACK | DEEPFOOL | FGSM |
|---|---|---|---|
| D3 ($\eta = 40K, \kappa = 5$) | **94.0%** | **92.1** | **93.3** |
| D3 ($\eta = 10K, \kappa = 5$) | 93.9 | 91.8 | 93.1 |
| D3 ($\eta = 10K, \kappa = 4$) | 93.2 | 91.8 | 92.2 |

Table 5: Top-1 classification accuracy (%) on grey-box attacks (50 ImageNet classes)

| NETWORKS | NO ATTACK | DEEPFOOL | FGSM |
|---|---|---|---|
| D3 ($\eta = 40K, \kappa = 5$) | **94.0** | 86.0 | 91.8 |
| D3 ($\eta = 10K, \kappa = 5$) | 93.9 | 88.0 | **91.8** |
| D3 ($\eta = 10K, \kappa = 4$) | 93.2 | **89.7** | 91.7 |

Table 6: Top-1 classification accuracy (%) on white-box attacks (50 ImageNet classes)

| NETWORKS | NO ATTACK | DEEPFOOL | FGSM |
|---|---|---|---|
| D3 ($\eta = 40K, \kappa = 5$) | **93.3** | 66.2 | 85.4 |
| D3 ($\eta = 10K, \kappa = 5$) | 92.8 | 66.4 | 85.4 |
| D3 ($\eta = 10K, \kappa = 4$) | 91.7 | **70.9** | **85.8** |

Table 7: Classification accuracy (%) CIFAR-10 dataset.

| NETWORK | TRAINING/DEFENSE | NO ATTACK | FGSM |
|---|---|---|---|
| RESNET | NORMAL | 92 | 11 |
| VGG | NORMAL | 89 | 30 |
| RESNET | ADVERSARIAL FGSM | 91 | 91 |
| RESNET | ADVERSARIAL BIM | 87 | 34 |
| RESNET | LABEL SMOOTHING | 92 | 28 |
| RESNET | FEATURE SQUEEZING | 84 | 18 |
| RESNET | ADVERSARIAL FGSM + FEATURE SQUEEZING | 86 | 55 |
| RESNET | NORMAL + PIXELDEFEND | 85 | 24 |
| VGG | NORMAL + PIXELDEFEND | 82 | 52 |
| RESNET | ADVERSARIAL FGSM + PIXELDEFEND | 88 | 67 |
| RESNET | ADVERSARIAL FGSM + ADAPTIVE PIXELDEFEND | 90 | 67 |
| RESNET | D3 | 87 | 80 |

## Appendix B: Experiments on CIFAR-10 dataset

Table 7 compares the D3 algorithm withe several recent defense algorithms under FGSM attack. Evaluation of the defense algorithms is taken from PixelDefend [45](PD) paper. The D3 performs better than all the methods except Adversarial FGSM which is specifically trained using FGSM perturbed images, while our method is attack-agnostic.