

# Appendix - Shape Matters: Deformable Patch Attack

Zhaoyu Chen<sup>1</sup>, Bo Li<sup>2†</sup>, Shuang Wu<sup>2</sup>, Jianghe Xu<sup>2</sup>,  
Shouhong Ding<sup>2</sup>, and Wenqiang Zhang<sup>1,3</sup>

<sup>1</sup> Academy for Engineering and Technology, Fudan University

<sup>2</sup> Youtu Lab, Tencent

<sup>3</sup> Yiwu Research Institute of Fudan University

## 1 Overview

This document provides more details of our **DAPatch**, organized as follows:

- In Section 2, we describe the whole algorithm of Deformable Adversarial Patch, corresponding to Section 3.3 of the main body.
- In Section 3, we supplement the description of the experimental setup, corresponding to Section 4.1 of the main body.
- In Section 4, we show the complete data of white box attacks under different areas, corresponding to Section 4.2, 4.3 and 4.6 of the main body.
- In Section 5, we add physical attack examples at different angles and lighting, corresponding to Section 4.5 of the main body.
- In Section 6, we show the ablation study on some hyper-parameters, including the sparsity of activation function, the shape loss  $L_{shape}$ , the number of rays  $R$ , and shape ratio  $s$ .
- In Section 7, we provide more visual comparison.

## 2 Deformable Adversarial Patch

Our proposed Deformable Adversarial Patch is summarized as Algorithm 1.

## 3 Experimental Setup

### 3.1 Comparable Methods

We use circular and square shape initialization in both GAP [1] and LaVAN [5]. We use random noise as initialization for both untargeted and targeted attacks. For all experiments, we set the number of iterations to 100 and  $\gamma$  to 1.  $u$  is the value of 1 pixel value after normalization,  $\alpha$  equals to  $8u$  before perturbation tuning, and  $\alpha$  equals to  $u$  after perturbation tuning. For PS-GAN [6], we use PS-GAN with weak constraints. The settings and constraints are the same as [6].

---

<sup>†</sup> indicates the corresponding author ([libraboli@tencent.com](mailto:libraboli@tencent.com)).

**Algorithm 1** Deformable Adversarial Patch (DAPatch)

**Input:** image  $x \in [x_{\min}, x_{\max}]$ , label  $y$ , the center  $O$ , the number of rays  $R$ , ray length array  $r = \{r_1, r_2, \dots, r_n\}$ , shape ratio  $s$ , patch percent  $pc$ , perturbation step  $\alpha$ , regular parameter  $\beta$ , ray step  $\gamma$ , the number of iteration  $T$

**Output:**  $x_{adv}^T$

---

```

1: Random sample  $\delta^0$  from  $[x_{\min}, x_{\max}]$ 
2:  $r^0 \leftarrow r, x_{adv}^0 \leftarrow x$ 
3: for  $k \in [1, T]$  do
4:   if  $k < \text{int}(s * T)$  then
5:      $M^k \leftarrow \text{DRP}(O, R, r^{k-1})$ 
6:      $x_{adv}^k \leftarrow (I - M^k) \odot x + M^k \odot \delta^{k-1}$ 
7:      $z^k \leftarrow f(x_{adv}^k)$ 
8:      $l \leftarrow L(z^k, y, pc, \beta)$ 
9:      $\delta^k \leftarrow \text{Clip}(\delta^{k-1} + \alpha \cdot \text{sign}(\nabla_{x_{adv}^k} l), x_{\min}, x_{\max})$ 
10:     $r^k \leftarrow \text{Clip}(r^{k-1} + \gamma \cdot \text{sign}(\nabla_{r^{k-1}} l), 1, \infty)$ 
11:    else if  $k == \text{int}(s * T)$  then
12:       $M \leftarrow \text{Sharpen}(M^{k-1})$ 
13:    else
14:       $x_{adv}^k \leftarrow (I - M) \odot x + M \odot \delta^{k-1}$ 
15:       $z^k \leftarrow f(x_{adv}^k)$ 
16:       $l \leftarrow L(z^k, y, pc, \beta)$ 
17:       $\delta^k \leftarrow \text{Clip}(\delta^{k-1} + \alpha \cdot \text{sign}(\nabla_{x_{adv}^k} l), x_{\min}, x_{\max})$ 
18:    end if
19:  end for
20: return  $x_{adv}^T$ 

```

---

For a image classifier  $f : x \rightarrow y$ , we denote the clean image as  $x \in R^{c \times h \times w}$  and the adversarial image as  $x_{adv}^k \in R^{c \times h \times w}$  at the  $k$ -th iteration. We also denote the corresponding label as  $y$  and the predicted label as  $\hat{y}$ . In Algorithm 1, the loss function of GAP is expressed as:

$$L = \begin{cases} CE(x_{adv}^k, y), & \text{untargeted attack} \\ -CE(x_{adv}^k, \hat{y}), & \text{targeted attack} \end{cases}, \quad (1)$$

and the loss function of LaVAN is expressed as:

$$L = \begin{cases} CE(x_{adv}^k, y) - CE(x_{adv}^k, y_s), & \text{untargeted attack} \\ CE(x_{adv}^k, y) - CE(x_{adv}^k, y_t), & \text{targeted attack} \end{cases}, \quad (2)$$

where  $y_s$  is the highest class other than class  $y$  and  $y_t$  is the pre-set target class.

### 3.2 Adversarial Training

Adversarial training is currently the most mainstream and effective method in adversarial defense. We choose the most efficient and powerful adversarial training method as the threat model.

**Fast-AT** Fast-AT [11] shows that adversarial training with the fast gradient sign method (FGSM), when combined with random initialization, is as effective as PGD-based training but has significantly lower cost. In the experiment, we choose Fast-AT ( $\epsilon = 4/255$ ) as the benchmark model of the attack.

**Feature Denoising** Feature Denoising [12] is the state-of-the-art defense against traditional perturbation-based adversarial attacks in a white-box setting, which contains blocks that denoise the features using non-local means or other filters. On ImageNet, under 10-iteration PGD white-box attacks, it achieves 55.7%. Even under extreme 2000-iteration PGD white-box attacks, it secures 42.6% accuracy. In the experiment, we choose its three open source models (Adv-ResNet-152, ResNet-152-Denoise and Resnext-101-Deniose) as the benchmark models for the attack.

## 4 More Experimental Results

We demonstrate the effectiveness of our proposed DAPatch on models of different architectures. We divide the model architecture into three categories: **Convolutional Neural Network** (VGG19 [9], Resnet-152 [3], DenseNet-161 [4] and MobileNet V2 [8]), **Vision Transformer** (ViT-B/16 [2] and Swin-B [7]), and **Neural Architecture Search** (EfficientNet-b7 [10]).

Table 2 illustrates that a particular shape can provide attack performance when textures are disabled. The area of the DAPatch shape is smaller than its convex hull, but it achieves a higher ASR. The experimental results in untargeted setting under 5 different patch areas on ILSVRC2012 and GTSRB are summarized in Table 3 and Table 4. For the more challenging targeted setting, the experimental results on ILSVRC2012 are reported in Table 5. The results of untargeted attacks on shape and texture bias are shown in Table 6. The results of untargeted attacks on adversarial training are shown in Table 7. All experiments show that when the patch area is small, DAPatch always obtains a higher ASR with a smaller area. Furthermore, under different patch areas, DAPatch can always obtain better attack performance within a smaller area compared with state-of-the-art methods.

## 5 Physical Attack

In this section, we provide more visual details about physical attacks. Figure 1 shows the different class example of DAPatch. Figure 2 shows the examples of DAPatch under different angles and lightning. So please **zoom** Figure 1 and Figure 2 to get more clearer shape details.



Fig. 1. More examples of physical attacks of DAPatch in untargeted setting.



(a) Middle,  $-30^\circ$     (b) Middle,  $0^\circ$     (c) Middle,  $30^\circ$     (d) Low,  $0^\circ$     (e) High,  $0^\circ$

Fig. 2. More examples of DAPatch with different angles and lighting.

## 6 Ablation Study

### 6.1 Ablation Study on $\lambda$

We review the special activation function  $\Phi$ , which is expressed as:

$$\Phi(x) = \frac{\tanh(\lambda(x-1)) + 1}{2}. \quad (3)$$

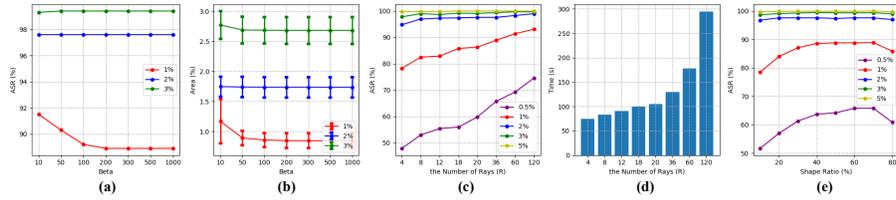
The  $\lambda$  controls the sparsity of activation function. Here, we study the effect of  $\lambda$  on attack performance in 2% area under the untargeted setting, as shown in Table 1. Experiments show that when  $\lambda = -100$ , the attack performance is better.  $\tanh$  can better make the mask close to binarization, so it is selected as the activation function.

Table 1. Ablation study on  $\lambda$ .

$\lambda$	-10	-50	-100	-300	-500
MoblieNet v2	89.9	94.9	<b>97.6</b>	93.6	92.6
Vit-B/16-224	75.5	88.9	<b>95.0</b>	86.1	84.8
ResNet-152	77.1	86.8	<b>93.1</b>	84.0	83.1

### 6.2 Ablation Study on Shape Loss $L_{shape}$

The area of patches needs to be controlled by  $L_{shape}$ . We compare the area variation wrt  $L_{shape}$  in Figure 3. We find that, if there is no  $L_{shape}$ , the area of the patch will increase indefinitely, but this does not meet the experimental settings. According to Figure 3, When  $\beta = 10$ ,  $L_{shape}$  cannot control the area well. When the  $\beta$  is large, the patch area can be better constrained to be within the specified percentage. Note that when  $\beta$  is large, ASR and area are not sensitive to beta, so in the experiment, to better control the area, we set  $\beta = 200$ .



**Fig. 3.** Ablation study. (a) shows the relationship between  $\beta$  and ASR. Error bars in (b) represent the standard deviation of the area of DAPatch. We can find when  $\beta$  is large, ASR and area are not sensitive to  $\beta$ . (c) and (d) are the ASR upper bound analysis. More rays can model patches with higher ASR and the result is saturated with 120 rays. The time cost increases as  $R$  becomes larger. (e) is the ablation study on the shape ratio  $s$ . It has the best attack performance at about  $s = 70$ .

### 6.3 Ablation Study on the Number of Rays $R$

It plays a fundamental role in the DAPatch and explicitly affects the shape modeling ability of the patch. From Figure 3 (c), more rays show higher upper bound and better ASR. For example, 36 rays improve by 6.0% ASR compared to 20 rays in 0.5% patch percentage. The 120 rays also saturate the performance since it depicts the patches well already and the rays are not the only constraint. Note that according to Figure 3 (d), the performance of 120 rays is not much improved compared to 36 rays, which is higher than other baseline methods, but it brings more than double the training time. In practice, considering for the efficiency, we set  $R = 36$ .

### 6.4 Ablation Study on Shape Ratio $s$

The shape ratio  $s$  is an important parameter to perturbation tune in DAPatch. Therefore, we evaluate the untargeted attack performance concerning different shape ratios  $s$  on MobileNet v2 in Figure 3 (e). When the patch area is small,  $s$  will greatly affect the attack effect. When the patch area is large, the effect of  $s$  is not very obvious. In practice, we choose  $s = 70$ .

## 7 More Visual Comparison

In this section, we provide more visual details. Figure 4 shows the visualization of DAPatch and other patch attacks. **Single** represents the single-anchor deformable patch representation and **Multi** represents the multi-anchor deformable patch representation. Figure 5 shows the deformation process of DAPatch in untargeted attacks under 5% area. **Disabling texture** means we just deform the shape and keep perturbations as white. The patches are generated on Mobilenet v2. The patches generated by multi-anchor deformable patch representation have more complex shapes, and there are cases where rays and contours intersect multiple times and the interior is hollowed out. So please **zoom** Figure 4 and Figure 5 to get more clearer shape details.

**Table 2.** The area of the convex hull is larger than DAPtach, but the attack performance is not as good as it, which shows that having a specific shape can improve the attack performance.

Network	Shape	0.5%		1%		2%		3%		5%	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
MoblieNet v2	Circle	1.5	0.510	2.2	0.964	4.5	2.040	6.8	3.031	10.2	4.982
	Square	1.4	0.504	1.7	1.054	3.3	2.010	4.4	3.023	5.9	4.888
	Ours	<b>8.9</b>	<b>0.377</b>	<b>13.4</b>	<b>0.790</b>	<b>21.0</b>	<b>1.648</b>	<b>25.8</b>	<b>2.496</b>	<b>35.7</b>	<b>4.340</b>
	Convex hull	2.4	0.902	4.0	2.143	6.4	4.090	9.1	5.676	12.4	8.296
Vit-B/16-224	Circle	0.9	0.510	1.4	0.964	2.2	2.040	2.2	3.031	2.7	4.982
	Square	0.5	0.504	0.7	1.054	1.1	2.010	1.6	3.023	2.0	4.888
	Ours	<b>8.6</b>	<b>0.355</b>	<b>12.0</b>	<b>0.789</b>	<b>16.3</b>	<b>1.563</b>	<b>20.7</b>	<b>2.507</b>	<b>27.2</b>	<b>4.247</b>
	Convex hull	1.3	0.720	2.4	1.874	3.4	3.644	4.4	5.149	5.5	7.643
ResNet-152	Circle	0.9	0.510	1.2	0.964	2.6	2.040	3.3	3.031	4.6	4.982
	Square	0.5	0.504	0.6	1.054	0.8	2.010	1.3	3.023	2.0	4.888
	Ours	<b>5.8</b>	<b>0.371</b>	<b>10.3</b>	<b>0.776</b>	<b>18.4</b>	<b>1.618</b>	<b>23.6</b>	<b>2.449</b>	<b>27.2</b>	<b>4.251</b>
	Convex hull	1.0	0.880	1.5	2.102	3.6	4.110	5.0	5.745	6.7	8.271

**Table 3.** Untargeted attacks of various network architectures on ILSVRC2012.

Network	Method	$\approx 0.5\%$		$\approx 1\%$		$\approx 2\%$		$\approx 3\%$		$\approx 5\%$	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
VGG-19	GAP <sub>s</sub>	73.4	0.510	92.6	0.964	98.4	2.040	99.3	3.031	99.7	4.982
	GAP <sub>c</sub>	72.4	0.504	94.5	1.054	98.7	2.010	99.3	3.023	99.8	4.888
	LaVAN <sub>s</sub>	76.9	0.510	92.6	0.964	98.9	2.040	99.5	3.031	100.0	4.982
	LaVAN <sub>c</sub>	78.5	0.504	95.1	1.054	99.0	2.010	99.5	3.023	100.0	4.888
	PS-GAN	74.5	0.510	94.2	0.964	97.4	2.040	99.2	3.031	100.0	4.982
	Ours	<b>78.6</b>	<b>0.449</b>	<b>95.6</b>	<b>0.868</b>	<b>99.1</b>	<b>1.744</b>	<b>99.5</b>	<b>2.759</b>	<b>100.0</b>	<b>4.598</b>
ResNet-152	GAP <sub>s</sub>	44.3	0.510	71.0	0.964	89.5	2.040	96.5	3.031	99.5	4.982
	GAP <sub>c</sub>	44.8	0.504	74.4	1.054	91.2	2.010	97.8	3.023	99.7	4.888
	LaVAN <sub>s</sub>	43.7	0.510	67.5	0.964	88.3	2.040	95.9	3.031	99.6	4.982
	LaVAN <sub>c</sub>	43.5	0.504	71.2	1.054	90.4	2.010	96.8	3.023	99.8	4.888
	PS-GAN	44.5	0.510	68.9	0.964	91.3	2.040	97.4	3.031	99.7	4.982
	Ours	<b>52.2</b>	<b>0.409</b>	<b>78.8</b>	<b>0.845</b>	<b>93.1</b>	<b>1.699</b>	<b>97.9</b>	<b>2.623</b>	<b>99.8</b>	<b>4.546</b>
DenseNet-161	GAP <sub>s</sub>	48.6	0.510	74.4	0.964	94.6	2.040	97.9	3.031	99.8	4.982
	GAP <sub>c</sub>	49.6	0.504	79.8	1.054	94.5	2.010	98.6	3.023	99.6	4.888
	LaVAN <sub>s</sub>	46.3	0.510	73.6	0.964	93.5	2.040	97.6	3.031	99.8	4.982
	LaVAN <sub>c</sub>	48.0	0.504	78.1	1.054	93.5	2.010	98.0	3.023	100.0	4.888
	PS-GAN	47.5	0.510	77.7	0.964	93.7	2.040	98.5	3.031	99.9	4.982
	Ours	<b>55.5</b>	<b>0.417</b>	<b>83.2</b>	<b>0.851</b>	<b>96.4</b>	<b>1.718</b>	<b>99.0</b>	<b>2.656</b>	<b>100.0</b>	<b>4.546</b>
MoblieNet v2	GAP <sub>s</sub>	57.9	0.510	83.6	0.964	96.2	2.040	98.6	3.031	99.9	4.982
	GAP <sub>c</sub>	57.8	0.504	86.7	1.054	97.2	2.010	99.2	3.023	100.0	4.888
	LaVAN <sub>s</sub>	56.6	0.510	81.8	0.964	95.6	2.040	98.9	3.031	99.9	4.982
	LaVAN <sub>c</sub>	58.0	0.504	84.3	1.054	96.6	2.010	98.7	3.023	99.8	4.888
	PS-GAN	54.5	0.510	84.2	0.964	95.5	2.040	99.3	3.031	99.9	4.982
	Ours	<b>65.8</b>	<b>0.423</b>	<b>88.9</b>	<b>0.847</b>	<b>97.6</b>	<b>1.735</b>	<b>99.4</b>	<b>2.684</b>	<b>100.0</b>	<b>4.578</b>
Efficientnet-b7	GAP <sub>s</sub>	43.3	0.510	63.5	0.964	85.5	2.040	91.5	3.031	97.2	4.982
	GAP <sub>c</sub>	42.5	0.504	68.8	1.054	88.0	2.010	94.4	3.023	97.5	4.888
	LaVAN <sub>s</sub>	42.3	0.510	64.7	0.964	89.5	2.040	95.9	3.031	98.3	4.982
	LaVAN <sub>c</sub>	41.2	0.504	69.2	1.054	89.2	2.010	95.9	3.023	98.1	4.888
	PS-GAN	40.9	0.510	65.8	0.964	89.3	2.040	95.2	3.031	97.9	4.982
	Ours	<b>45.7</b>	<b>0.442</b>	<b>71.1</b>	<b>0.956</b>	<b>89.6</b>	<b>2.003</b>	<b>95.9</b>	<b>3.014</b>	<b>98.3</b>	<b>4.824</b>
Vit-B/16-224	GAP <sub>s</sub>	47.0	0.510	72.0	0.964	92.4	2.040	97.2	3.031	99.8	4.982
	GAP <sub>c</sub>	45.6	0.504	77.2	1.054	93.0	2.010	97.8	3.023	99.7	4.888
	LaVAN <sub>s</sub>	44.8	0.510	71.8	0.964	93.5	2.040	98.3	3.031	99.9	4.982
	LaVAN <sub>c</sub>	46.9	0.504	74.9	1.054	93.5	2.010	98.3	3.023	99.9	4.888
	PS-GAN	45.9	0.510	71.9	0.964	90.2	2.040	97.4	3.031	99.8	4.982
	Ours	<b>56.9</b>	<b>0.417</b>	<b>80.9</b>	<b>0.849</b>	<b>95.0</b>	<b>1.717</b>	<b>98.3</b>	<b>2.676</b>	99.9	<b>4.528</b>
Swin-B-224	GAP <sub>s</sub>	32.4	0.510	68.2	0.964	91.8	2.040	97.6	3.031	99.7	4.982
	GAP <sub>c</sub>	34.4	0.504	77.7	1.054	94.5	2.010	98.6	3.023	99.6	4.888
	LaVAN <sub>s</sub>	35.3	0.510	68.6	0.964	95.8	2.040	99.0	3.031	99.9	4.982
	LaVAN <sub>c</sub>	37.2	0.504	75.8	1.054	96.6	2.010	99.4	3.023	100.0	4.888
	PS-GAN	31.2	0.510	66.4	0.964	91.2	2.040	97.6	3.031	99.7	4.982
	Ours	<b>39.7</b>	<b>0.395</b>	<b>79.6</b>	<b>0.818</b>	<b>97.1</b>	<b>1.668</b>	<b>99.5</b>	<b>2.545</b>	<b>100.0</b>	<b>4.354</b>

**Table 4.** Untargeted attacks of various network architectures on GTSRB.

Network	Method	$\approx 0.5\%$		$\approx 1\%$		$\approx 2\%$		$\approx 3\%$		$\approx 5\%$	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
ResNet-152	GAP_s	13.6	0.510	21.4	0.964	37.8	2.040	56.0	3.031	76.0	4.982
	GAP_c	13.0	0.504	22.6	1.054	38.0	2.010	57.6	3.023	80.0	4.888
	LaVAN_s	14.6	0.510	22.8	0.964	41.6	2.040	60.4	3.031	83.8	4.982
	LaVAN_c	14.8	0.504	26.0	1.054	41.8	2.010	58.4	3.023	83.2	4.888
	PS-GAN	13.7	0.510	23.4	0.964	39.4	2.040	59.5	3.031	82.3	4.982
	Ours	<b>15.0</b>	<b>0.477</b>	<b>27.1</b>	<b>0.831</b>	<b>42.3</b>	<b>1.932</b>	<b>61.5</b>	<b>2.873</b>	<b>85.6</b>	<b>4.722</b>
Efficientnet-b7	GAP_s	20.4	0.510	45.0	0.964	68.0	2.040	84.0	3.031	94.4	4.982
	GAP_c	20.6	0.504	39.4	1.054	66.6	2.010	82.6	3.023	94.4	4.888
	LaVAN_s	22.2	0.510	46.2	0.964	74.0	2.040	89.2	3.031	95.6	4.982
	LaVAN_c	22.8	0.504	51.4	1.054	74.0	2.010	89.0	3.023	97.2	4.888
	PS-GAN	21.5	0.510	46.2	0.964	71.2	2.040	85.6	3.031	96.2	4.982
	Ours	<b>23.6</b>	<b>0.469</b>	<b>53.1</b>	<b>0.893</b>	<b>75.2</b>	<b>1.873</b>	<b>89.5</b>	<b>2.934</b>	<b>98.7</b>	<b>4.825</b>
Vit-B/16-224	GAP_s	28.6	0.510	61.2	0.964	90.0	2.040	97.6	3.031	99.8	4.982
	GAP_c	28.4	0.504	68.0	1.054	90.2	2.010	98.2	3.023	99.0	4.888
	LaVAN_s	28.2	0.510	61.6	0.964	91.8	2.040	98.2	3.031	99.8	4.982
	LaVAN_c	26.2	0.504	65.4	1.054	92.0	2.010	97.4	3.023	99.6	4.888
	PS-GAN	27.4	0.510	64.2	0.964	90.2	2.040	98.1	3.031	99.7	4.982
	Ours	<b>30.1</b>	<b>0.483</b>	<b>68.1</b>	<b>0.896</b>	<b>93.5</b>	<b>1.783</b>	<b>98.9</b>	<b>2.892</b>	<b>100.0</b>	<b>4.732</b>

**Table 5.** Targeted attacks of various network architectures on ILSVRC2012.

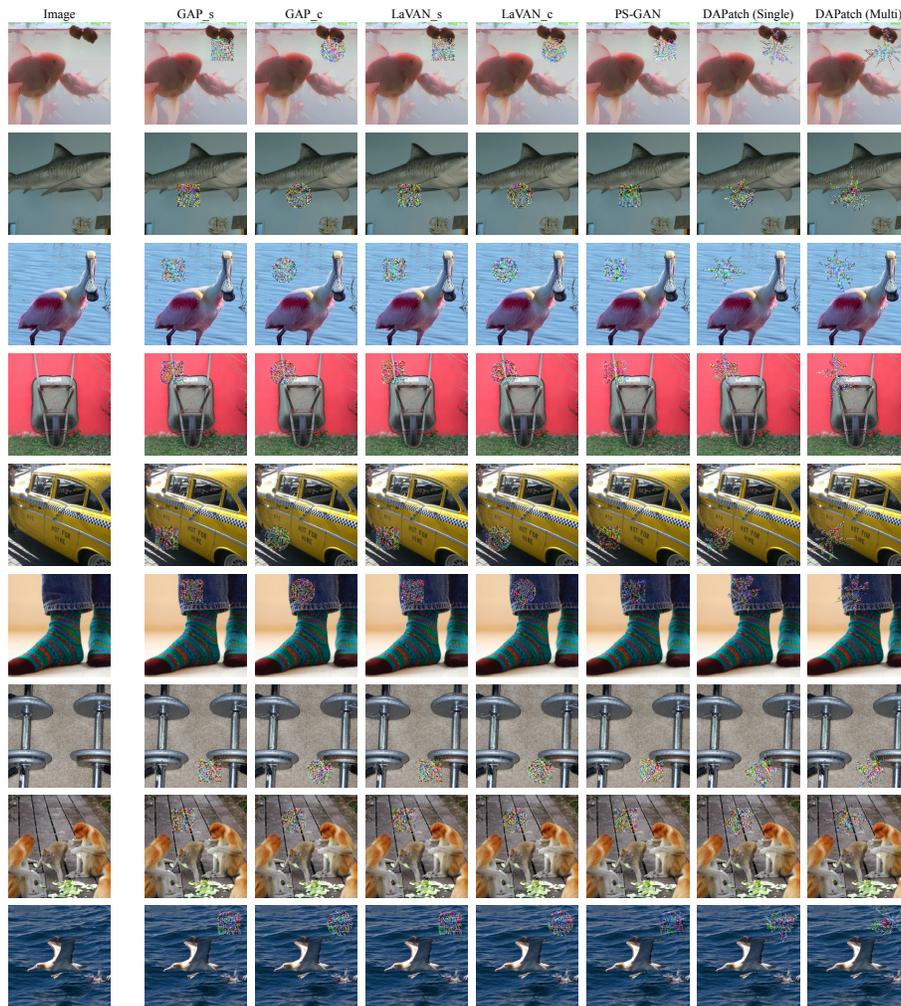
Network	Method	$\approx 1\%$		$\approx 3\%$		$\approx 5\%$		$\approx 7\%$		$\approx 10\%$	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
VGG-19	GAP <sub>s</sub>	16.60	0.964	59.1	3.031	81.0	4.982	93.9	6.938	96.7	10.046
	GAP <sub>c</sub>	21.40	1.054	58.5	3.023	77.3	4.888	90.0	6.794	98.1	10.002
	LaVAN <sub>s</sub>	5.60	0.964	24.3	3.031	37.4	4.982	52.0	6.938	63.0	10.046
	LaVAN <sub>c</sub>	5.20	1.054	22.8	3.023	36.3	4.888	47.0	6.794	58.4	10.002
	PS-GAN	20.60	0.964	59.4	3.031	81.4	4.982	94.2	6.938	97.2	10.046
	Ours	<b>21.40</b>	<b>0.864</b>	<b>61.7</b>	<b>2.676</b>	<b>82.8</b>	<b>4.560</b>	<b>94.6</b>	<b>6.521</b>	<b>98.3</b>	<b>9.225</b>
ResNet-152	GAP <sub>s</sub>	5.30	0.964	41.1	3.031	68.8	4.982	87.3	6.938	95.0	10.046
	GAP <sub>c</sub>	8.80	1.054	44.3	3.023	72.8	4.888	88.3	6.794	97.4	10.002
	LaVAN <sub>s</sub>	3.50	0.964	22.6	3.031	45.1	4.982	61.8	6.938	78.0	10.046
	LaVAN <sub>c</sub>	6.00	1.054	23.9	3.023	47.7	4.888	64.7	6.794	81.1	10.002
	PS-GAN	8.90	0.964	46.2	3.031	73.8	4.982	87.4	6.938	96.4	10.046
	Ours	<b>9.10</b>	<b>0.849</b>	<b>48.7</b>	<b>2.668</b>	<b>78.1</b>	<b>4.553</b>	<b>90.2</b>	<b>6.439</b>	<b>97.6</b>	<b>9.232</b>
DenseNet-161	GAP <sub>s</sub>	9.20	0.964	45.0	3.031	73.3	4.982	92.2	6.938	96.9	10.046
	GAP <sub>c</sub>	11.50	1.054	44.9	3.023	73.1	4.888	90.4	6.794	98.3	10.002
	LaVAN <sub>s</sub>	4.60	0.964	21.6	3.031	35.9	4.982	51.0	6.938	67.9	10.046
	LaVAN <sub>c</sub>	5.90	1.054	24.1	3.023	38.2	4.888	54.0	6.794	72.4	10.002
	PS-GAN	10.30	0.964	49.2	3.031	74.5	4.982	92.3	6.938	97.5	10.046
	Ours	<b>13.70</b>	<b>0.858</b>	<b>55.3</b>	<b>2.682</b>	<b>77.8</b>	<b>4.569</b>	<b>92.6</b>	<b>6.472</b>	<b>98.6</b>	<b>9.315</b>
MobileNet v2	GAP <sub>s</sub>	4.50	0.964	51.3	3.031	81.1	4.982	93.7	6.938	97.3	10.046
	GAP <sub>c</sub>	5.80	1.054	51.6	3.023	81.7	4.888	94.8	6.794	99.3	10.002
	LaVAN <sub>s</sub>	1.60	0.964	30.0	3.031	51.7	4.982	71.8	6.938	86.2	10.046
	LaVAN <sub>c</sub>	2.60	1.054	31.8	3.023	52.8	4.888	72.5	6.794	88.1	10.002
	PS-GAN	5.90	0.964	52.4	3.031	82.3	4.982	93.9	6.938	99.5	10.046
	Ours	<b>7.60</b>	<b>0.869</b>	<b>54.8</b>	<b>2.271</b>	<b>84.9</b>	<b>4.593</b>	<b>94.8</b>	<b>6.481</b>	<b>99.7</b>	<b>9.433</b>
Efficientnet-b7	GAP <sub>s</sub>	4.40	0.964	52.1	3.031	81.9	4.982	93.7	6.938	97.2	10.046
	GAP <sub>c</sub>	6.20	1.054	53.6	3.023	81.4	4.888	93.4	6.794	99.0	10.002
	LaVAN <sub>s</sub>	1.50	0.964	33.1	3.031	65.0	4.982	82.2	6.938	91.6	10.046
	LaVAN <sub>c</sub>	2.30	1.054	34.0	3.023	62.1	4.888	83.6	6.794	95.1	10.002
	PS-GAN	4.90	0.964	53.6	3.031	81.5	4.982	93.2	6.938	99.1	10.046
	Ours	<b>7.60</b>	<b>0.869</b>	<b>53.6</b>	<b>2.953</b>	<b>82.0</b>	<b>4.851</b>	<b>93.7</b>	<b>6.713</b>	<b>99.3</b>	<b>10.000</b>
Vit-B/16-224	GAP <sub>s</sub>	6.20	0.964	48.8	3.031	85.4	4.982	97.3	6.938	97.9	10.046
	GAP <sub>c</sub>	7.90	1.054	50.6	3.023	85.7	4.888	97.2	6.794	100.0	10.002
	LaVAN <sub>s</sub>	3.10	0.964	25.4	3.031	52.8	4.982	78.3	6.938	93.4	10.046
	LaVAN <sub>c</sub>	4.20	1.054	24.7	3.023	54.9	4.888	78.4	6.794	96.6	10.002
	PS-GAN	5.30	0.964	49.3	3.031	85.8	4.982	97.1	6.938	98.1	10.046
	Ours	<b>9.60</b>	<b>0.850</b>	<b>50.7</b>	<b>2.697</b>	<b>86.2</b>	<b>4.688</b>	<b>97.4</b>	<b>6.727</b>	<b>100.0</b>	<b>9.322</b>
Swin-B-224	GAP <sub>s</sub>	18.50	0.964	91.2	3.031	99.5	4.982	100.0	6.938	100.0	10.046
	GAP <sub>c</sub>	27.80	1.054	94.4	3.023	99.5	4.888	99.9	6.794	100.0	10.002
	LaVAN <sub>s</sub>	19.00	0.964	85.1	3.031	97.0	4.982	99.4	6.938	98.0	10.046
	LaVAN <sub>c</sub>	25.20	1.054	86.5	3.023	96.9	4.888	99.3	6.794	100.0	10.002
	PS-GAN	29.60	0.964	95.1	3.031	98.9	4.982	99.4	6.938	100.0	10.046
	Ours	<b>38.20</b>	<b>0.848</b>	<b>98.5</b>	<b>2.587</b>	<b>99.7</b>	<b>4.371</b>	<b>100.0</b>	<b>6.216</b>	<b>100.0</b>	<b>8.912</b>

**Table 6.** Untargeted attacks of shape and texture bias on ILSVRC2012.

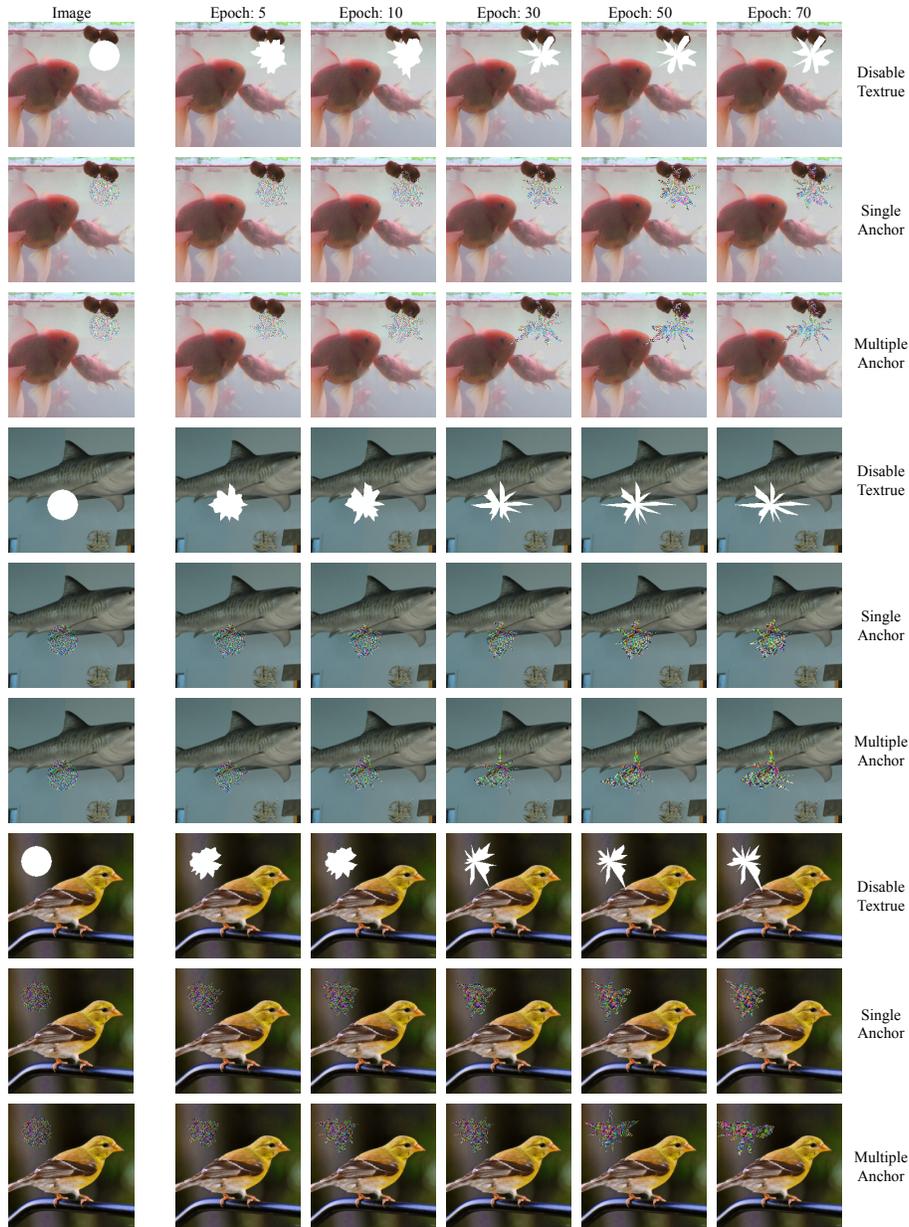
Network	Method	$\approx 0.5\%$		$\approx 1\%$		$\approx 2\%$		$\approx 3\%$		$\approx 5\%$	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
ResNet50-SIN	GAP_s	70.4	0.510	87.3	0.964	96.9	2.040	99.1	3.031	99.8	4.882
	GAP_c	70.1	0.504	88.5	1.054	96.9	2.010	99.6	3.023	99.8	4.888
	LaVAN_s	66.2	0.510	82.2	0.964	95.1	2.040	98.1	3.031	99.8	4.982
	LaVAN_c	65.6	0.504	84.2	1.054	96.0	2.010	98.7	3.023	99.7	4.888
	PS-GAN	70.0	0.510	85.3	0.964	96.8	2.040	99.5	3.031	99.8	4.982
	Ours	<b>74.1</b>	<b>0.446</b>	<b>90.3</b>	<b>0.893</b>	<b>98.7</b>	<b>1.764</b>	<b>99.6</b>	<b>2.724</b>	<b>99.9</b>	<b>4.620</b>
ResNet50-SIN+IN	GAP_s	44.3	0.510	68.3	0.964	90.6	2.040	95.7	3.031	98.5	4.982
	GAP_c	44.7	0.504	72.4	1.054	90.9	2.010	96.3	3.023	99.1	4.888
	LaVAN_s	41.2	0.510	64.9	0.964	88.8	2.040	94.6	3.031	98.7	4.982
	LaVAN_c	42.4	0.504	69.4	1.054	89.4	2.010	96.1	3.023	98.9	4.888
	PS-GAN	44.5	0.510	68.9	0.964	90.6	2.040	95.2	3.031	98.2	4.982
	Ours	<b>48.1</b>	<b>0.426</b>	<b>75.6</b>	<b>0.860</b>	<b>91.5</b>	<b>1.750</b>	<b>96.3</b>	<b>2.669</b>	<b>99.1</b>	<b>4.587</b>
ResNet50-SIN+IN-IN	GAP_s	41.6	0.510	62.2	0.964	85.6	2.040	93.0	3.031	98.0	4.982
	GAP_c	44.3	0.504	68.6	1.054	87.4	2.010	94.2	3.023	98.3	4.888
	LaVAN_s	38.7	0.510	58.2	0.964	83.1	2.040	92.8	3.031	98.2	4.982
	LaVAN_c	39.6	0.504	65.3	1.054	84.7	2.010	93.4	3.023	98.2	4.888
	PS-GAN	39.2	0.510	62.7	0.964	85.2	2.040	94.7	3.031	98.6	4.982
	Ours	<b>44.3</b>	<b>0.420</b>	<b>70.1</b>	<b>0.845</b>	<b>88.9</b>	<b>1.729</b>	<b>95.1</b>	<b>2.651</b>	<b>99.0</b>	<b>4.562</b>
ResNet50-Debiased	GAP_s	42.5	0.510	64.2	0.964	85.8	2.040	93.1	3.031	98.4	4.982
	GAP_c	44.0	0.504	68.6	1.054	87.9	2.010	94.3	3.023	98.4	4.888
	LaVAN_s	38.2	0.510	59.7	0.964	83.3	2.040	90.9	3.031	97.9	4.982
	LaVAN_c	38.8	0.504	64.2	1.054	84.3	2.010	93.2	3.023	98.6	4.888
	PS-GAN	43.2	0.510	67.6	0.964	88.0	2.040	93.4	3.031	98.0	4.982
	Ours	<b>46.6</b>	<b>0.438</b>	<b>68.6</b>	<b>0.852</b>	<b>88.1</b>	<b>1.744</b>	<b>94.7</b>	<b>2.665</b>	<b>98.6</b>	<b>4.593</b>
ResNet152-Debiased	GAP_s	33.3	0.510	53.0	0.964	81.4	2.040	91.0	3.031	98.3	4.982
	GAP_c	34.1	0.504	58.9	1.054	84.1	2.010	93.0	3.023	98.4	4.888
	LaVAN_s	30.6	0.510	49.8	0.964	77.0	2.040	88.8	3.031	98.1	4.982
	LaVAN_c	29.8	0.504	52.1	1.054	77.5	2.010	90.1	3.023	98.4	4.888
	PS-GAN	32.3	0.510	53.4	0.964	83.0	2.040	93.2	3.031	98.3	4.982
	Ours	<b>37.4</b>	<b>0.422</b>	<b>61.8</b>	<b>0.850</b>	<b>85.3</b>	<b>1.735</b>	<b>94.5</b>	<b>2.626</b>	<b>98.5</b>	<b>4.532</b>

**Table 7.** Untargeted attacks on networks with adversarial training.

Network	Method	$\approx 0.5\%$		$\approx 1\%$		$\approx 2\%$		$\approx 3\%$		$\approx 5\%$	
		ASR	Area	ASR	Area	ASR	Area	ASR	Area	ASR	Area
Adv-ResNet-152	GAP_s	61.0	0.510	74.5	0.964	86.6	2.040	90.3	3.031	94.4	4.982
	GAP_c	60.6	0.504	77.4	1.054	87.2	2.010	91.7	3.023	95.1	4.888
	LaVAN_s	58.4	0.510	71.1	0.964	83.9	2.040	88.8	3.031	93.8	4.982
	LaVAN_c	57.2	0.504	72.6	1.054	83.9	2.010	89.3	3.023	94.8	4.888
	PS-GAN	59.2	0.510	77.7	0.964	84.3	2.040	89.7	3.031	95.2	4.982
	Ours	<b>62.5</b>	<b>0.472</b>	<b>78.4</b>	<b>0.948</b>	<b>88.2</b>	<b>1.921</b>	<b>92.4</b>	<b>2.791</b>	<b>96.6</b>	<b>4.703</b>
ResNet-152-Denoise	GAP_s	59.3	0.510	74.5	0.964	86.5	2.040	92.6	3.031	96.4	4.982
	GAP_c	59.0	0.504	77.3	1.054	87.8	2.010	92.9	3.023	96.4	4.888
	LaVAN_s	59.6	0.510	72.6	0.964	84.7	2.040	91.8	3.031	96.9	4.982
	LaVAN_c	60.7	0.504	75.1	1.054	85.5	2.010	92.7	3.023	96.5	4.888
	PS-GAN	61.7	0.510	75.1	0.964	86.2	2.040	92.2	3.031	96.3	4.982
	Ours	<b>62.3</b>	<b>0.464</b>	<b>77.4</b>	<b>0.959</b>	<b>88.0</b>	<b>1.835</b>	<b>92.9</b>	<b>2.853</b>	<b>98.3</b>	<b>4.619</b>
Resnext-101-Deniose	GAP_s	50.4	0.510	66.3	0.964	83.8	2.040	90.2	3.031	95.0	4.982
	GAP_c	51.1	0.504	70.5	1.054	84.2	2.010	89.9	3.023	95.7	4.888
	LaVAN_s	49.7	0.510	65.0	0.964	80.6	2.040	87.6	3.031	94.9	4.982
	LaVAN_c	49.5	0.504	67.9	1.054	81.2	2.010	88.4	3.023	95.2	4.888
	PS-GAN	51.2	0.510	68.1	0.964	80.9	2.040	89.9	3.031	94.7	4.982
	Ours	<b>52.9</b>	<b>0.471</b>	<b>68.9</b>	<b>0.949</b>	<b>85.5</b>	<b>1.928</b>	<b>90.2</b>	<b>2.814</b>	<b>95.7</b>	<b>4.706</b>
Fast_AT	GAP_s	50.4	0.510	62.3	0.964	80.4	2.040	88.5	3.031	95.0	4.982
	GAP_c	50.6	0.504	65.5	1.054	80.3	2.010	88.8	3.023	94.9	4.888
	LaVAN_s	48.7	0.510	60.2	0.964	77.5	2.040	84.7	3.031	92.9	4.982
	LaVAN_c	48.7	0.504	62.6	1.054	78.0	2.010	85.3	3.023	93.1	4.888
	PS-GAN	48.9	0.510	63.4	0.964	79.1	2.040	85.2	3.031	93.2	4.982
	Ours	<b>51.3</b>	<b>0.473</b>	<b>65.6</b>	<b>0.944</b>	<b>82.0</b>	<b>1.890</b>	<b>90.0</b>	<b>2.963</b>	<b>95.4</b>	<b>4.772</b>



**Fig. 4.** Visualization of DAPatch and other patch attacks under 5% area. Please **zoom** images for better shape details.



**Fig. 5.** Deformation process of DAPatch in untargeted attacks under 5% area. Please **zoom** images for better shape details.

## References

1. Brown, T.B., Mané, D., Roy, A., Abadi, M., Gilmer, J.: Adversarial patch (2017), <http://arxiv.org/abs/1712.09665> 1
2. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net (2021), <https://openreview.net/forum?id=YicbFdNTTy> 3
3. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016. pp. 770–778. IEEE Computer Society (2016). <https://doi.org/10.1109/CVPR.2016.90>, <https://doi.org/10.1109/CVPR.2016.90> 3
4. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017. pp. 2261–2269. IEEE Computer Society (2017). <https://doi.org/10.1109/CVPR.2017.243>, <https://doi.org/10.1109/CVPR.2017.243> 3
5. Karmon, D., Zoran, D., Goldberg, Y.: Lavan: Localized and visible adversarial noise. In: Dy, J.G., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018. Proceedings of Machine Learning Research, vol. 80, pp. 2512–2520. PMLR (2018), <http://proceedings.mlr.press/v80/karmon18a.html> 1
6. Liu, A., Liu, X., Fan, J., Ma, Y., Zhang, A., Xie, H., Tao, D.: Perceptual-sensitive GAN for generating adversarial patches. In: The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019. pp. 1028–1035. AAAI Press (2019). <https://doi.org/10.1609/aaai.v33i01.33011028>, <https://doi.org/10.1609/aaai.v33i01.33011028> 1
7. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows (2021), <https://arxiv.org/abs/2103.14030> 3
8. Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A., Chen, L.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 4510–4520. IEEE Computer Society (2018). <https://doi.org/10.1109/CVPR.2018.00474>, [http://openaccess.thecvf.com/content\\_cvpr\\_2018/html/Sandler\\_MobileNetV2\\_Inverted\\_Residuals\\_CVPR\\_2018\\_paper.html](http://openaccess.thecvf.com/content_cvpr_2018/html/Sandler_MobileNetV2_Inverted_Residuals_CVPR_2018_paper.html) 3
9. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: Bengio, Y., LeCun, Y. (eds.) 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015), <http://arxiv.org/abs/1409.1556> 3
10. Tan, M., Le, Q.V.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: Chaudhuri, K., Salakhutdinov, R. (eds.) Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long

- Beach, California, USA. Proceedings of Machine Learning Research, vol. 97, pp. 6105–6114. PMLR (2019), <http://proceedings.mlr.press/v97/tan19a.html> 3
11. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), <https://openreview.net/forum?id=BJx040EFvH> 3
  12. Xie, C., Wu, Y., van der Maaten, L., Yuille, A.L., He, K.: Feature denoising for improving adversarial robustness. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019. pp. 501–509. Computer Vision Foundation / IEEE (2019). <https://doi.org/10.1109/CVPR.2019.00059>, [http://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Xie\\_Feature\\_Denoising\\_for\\_Improving\\_Adversarial\\_Robustness\\_CVPR\\_2019\\_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Xie_Feature_Denoising_for_Improving_Adversarial_Robustness_CVPR_2019_paper.html) 3