

## Supplementary - Deflecting Adversarial Attacks with Pixel Deflection

### Results on various classifiers

Original classification accuracy of each classifier on selected 1000 images is reported in the table. However, we omit the images that were originally incorrectly classified, thus the accuracy of clean images without defense is always 100%. Weights for each classifier were obtained from Tensorflow GitHub repository <sup>1</sup>.

Model	$ L_2 $	No Defense	With Defense	
			Single	Ens-10
ResNet-50, original classification 76%				
Clean	0.00	100	98.3	<b>98.9</b>
FGSM	0.05	20.0	79.9	<b>81.5</b>
IGSM	0.03	14.1	83.7	<b>83.7</b>
DFool	0.02	26.3	86.3	<b>90.3</b>
JSMA	0.02	25.5	91.5	<b>97.0</b>
LBFGS	0.02	12.1	88.0	<b>91.6</b>
C&W	0.04	04.8	92.7	<b>98.0</b>
VGG-19, original classification 71%				
Clean	0.00	100	99.8	<b>99.8</b>
FGSM	0.05	12.2	79.3	<b>81.3</b>
IGSM	0.04	9.79	79.2	<b>81.6</b>
DFool	0.01	23.7	83.9	<b>91.6</b>
JSMA	0.01	29.1	95.8	<b>98.5</b>
LBFGS	0.03	13.8	83.0	<b>93.9</b>
C&W	0.04	0.00	93.1	<b>97.6</b>
Inception-v3, original classification 78%				
Clean	0.00	100	98.1	<b>98.5</b>
FGSM	0.05	22.1	85.8	<b>87.1</b>
IGSM	0.04	15.5	<b>89.7</b>	89.1
DFool	0.02	27.2	82.6	<b>85.3</b>
JSMA	0.02	24.2	93.7	<b>98.6</b>
LBFGS	0.02	12.5	87.1	<b>91.0</b>
C&W	0.04	07.1	93.9	<b>98.5</b>

Table 1: Params:  $\sigma = 0.04$ , Window=10, Deflections=100

Top-1 accuracy on applying pixel deflection and wavelet denoising across various attack models.

<sup>1</sup><https://github.com/tensorflow/models/tree/master/research/slim#Pretrained>

## Comparison of small and large perturbations

Model	$ L_2 $	No Defense	With Defense	
			Single	Ens-10
<b>Clean</b>	0.00	100	98.3	<b>98.9</b>
<b>FGSM</b>	0.05	20.0	79.9	<b>81.5</b>
<b>IGSM</b>	0.03	14.1	83.7	<b>83.7</b>
<b>DFool</b>	0.02	26.3	86.3	<b>90.3</b>
<b>JSMA</b>	0.02	25.5	91.5	<b>97.0</b>
<b>LBFGS</b>	0.02	12.1	88.0	<b>91.6</b>
<b>C&amp;W</b>	0.04	04.8	92.7	<b>98.0</b>
Large perturbations				
<b>FGSM</b>	0.12	11.1	61.5	<b>70.4</b>
<b>IGSM</b>	0.09	11.1	62.5	<b>72.5</b>
<b>DFool</b>	0.08	08.0	82.4	<b>88.9</b>
<b>JSMA</b>	0.05	22.1	88.9	<b>92.1</b>
<b>LBFGS</b>	0.04	12.1	77.0	<b>89.0</b>

Table 2: Params:  $\sigma = 0.04$ , Window=10, Deflections=100

Top-1 accuracy on applying pixel deflection and wavelet denoising across various attack models. We evaluate non-efficient attacks at larger  $|L_P|$  which leave visible perturbations to show the robustness of our model.

## Comparison of various shrinkage

Model	Hard	VISU	SURE	Bayes
<b>Clean</b>	39.5	96.1	92.1	<b>98.9</b>
<b>FGSM</b>	35.9	63.8	79.7	<b>81.5</b>
<b>IGSM</b>	42.5	67.8	81.1	<b>83.7</b>
<b>DFool</b>	37.2	78.4	87.7	<b>90.3</b>
<b>JSMA</b>	39.9	93.0	93.0	<b>97.0</b>
<b>LBFGS</b>	37.2	81.1	90.4	<b>91.6</b>
<b>C&amp;W</b>	36.8	93.4	92.8	<b>98.0</b>

Table 3: Params:  $\sigma = 0.04$ , Window=10, Deflections=100

Comparison of various thresholding techniques, after application of pixel deflection.

In Table 3 we present a comparison of various shrinkage methods on wavelet coefficients after pixel deflection. All the results reported are for applying the given thresholding after pixel deflection. BayesShrink, which learns separate Gaussian parameters for each coefficient, does better than other soft-thresholding techniques. VisuShrink is a faster technique as it uses a universal threshold but that limits its applicability on some images. SUREShrink has been shown to perform well with compression but as evident, in our results, it is less well suited to denoising.

## Ablation studies of various parameters

Attack	$ L_2$	No Defense	With Defense			
		Window=10, Deflections $\longrightarrow$	<b>10</b>	<b>100</b>	<b>1K</b>	<b>10K</b>
<b>Clean</b>	0.00	100	<b>98.4</b>	98.1	94.7	80.3
<b>FGSM</b>	0.04	19.2	75.7	<b>79.7</b>	71.7	69.1
<b>IGSM</b>	0.03	13.8	78.4	<b>81.7</b>	75.2	71.2
<b>DFool</b>	0.02	25.0	83.7	<b>87.7</b>	81.0	77.0
<b>JSMA</b>	0.02	25.9	91.7	<b>93.0</b>	87.7	67.7
<b>LBFGS</b>	0.02	11.6	85.0	<b>90.3</b>	82.4	73.0
<b>C&amp;W</b>	0.04	05.2	89.4	<b>93.1</b>	86.8	69.7

Table 4: Top-1 accuracy with different deflections.

Attack	L2	No Defense	With Defense			
		Deflections=100, Window $\longrightarrow$	<b>5</b>	<b>10</b>	<b>50</b>	<b>100</b>
<b>Clean</b>	0.00	100	<b>98.6</b>	98.1	96.4	94.4
<b>FGSM</b>	0.04	19.2	<b>79.7</b>	<b>79.7</b>	78.4	76.7
<b>IGSM</b>	0.03	13.8	81.0	<b>81.7</b>	79.7	78.4
<b>DFool</b>	0.02	25.0	86.4	<b>87.7</b>	87.7	85.0
<b>JSMA</b>	0.02	25.9	92.3	<b>93.0</b>	91.7	90.3
<b>LBFGS</b>	0.02	11.6	89.4	<b>90.3</b>	89.0	88.1
<b>C&amp;W</b>	0.04	05.2	91.8	<b>93.1</b>	90.5	89.2

Table 5: Top-1 accuracy with different window sizes.

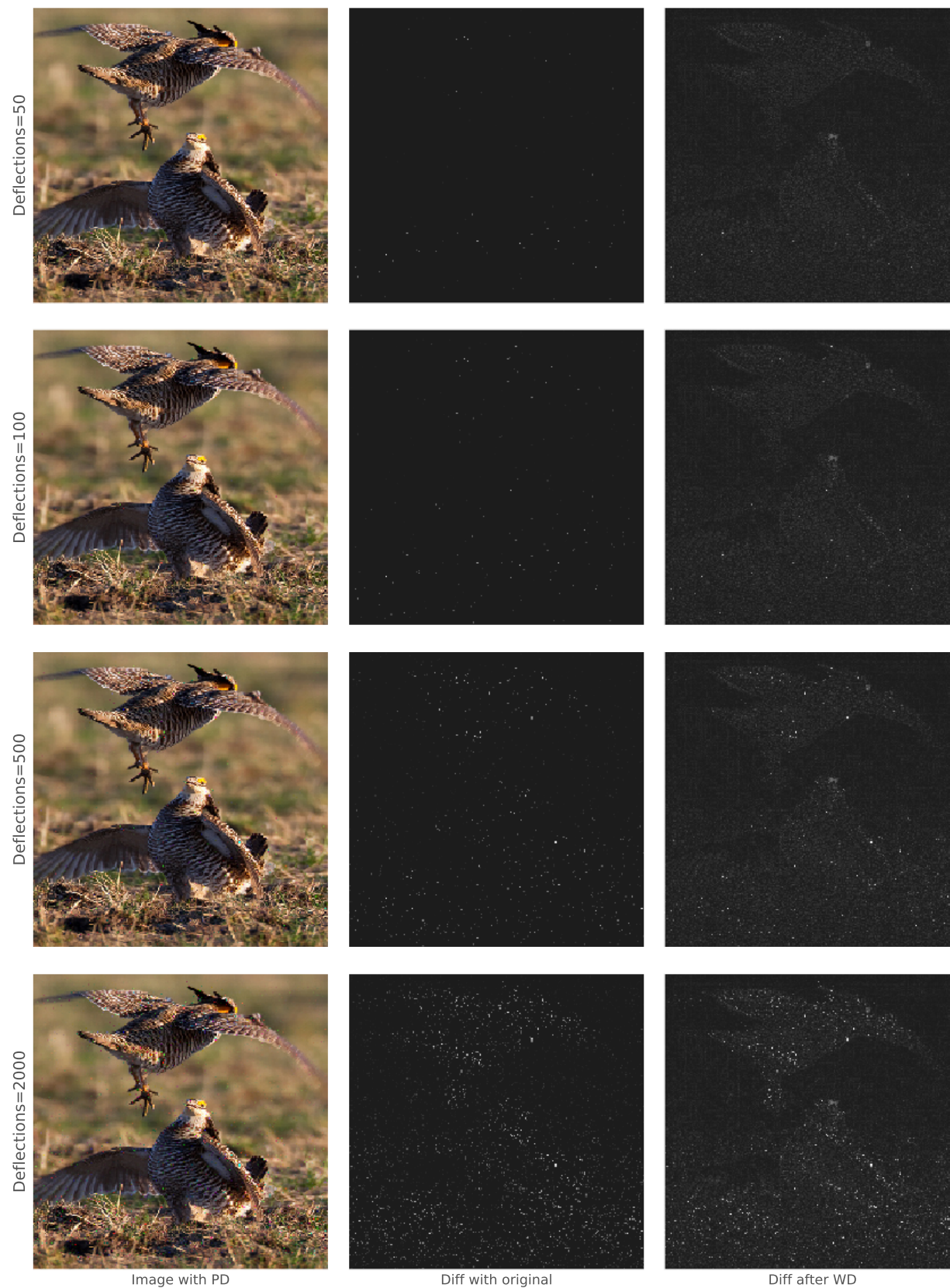


Figure 1: This is same image as Figure 1 on the paper but with a better color scheme. Impact of Pixel Deflection (PD) on a natural image and subsequent denoising using wavelet transform (WD). Left: Image with given number of pixels deflected. Middle: Difference between clean image and deflected image. Right: Difference between clean image and deflected image after denoising. Enlarge to see details.

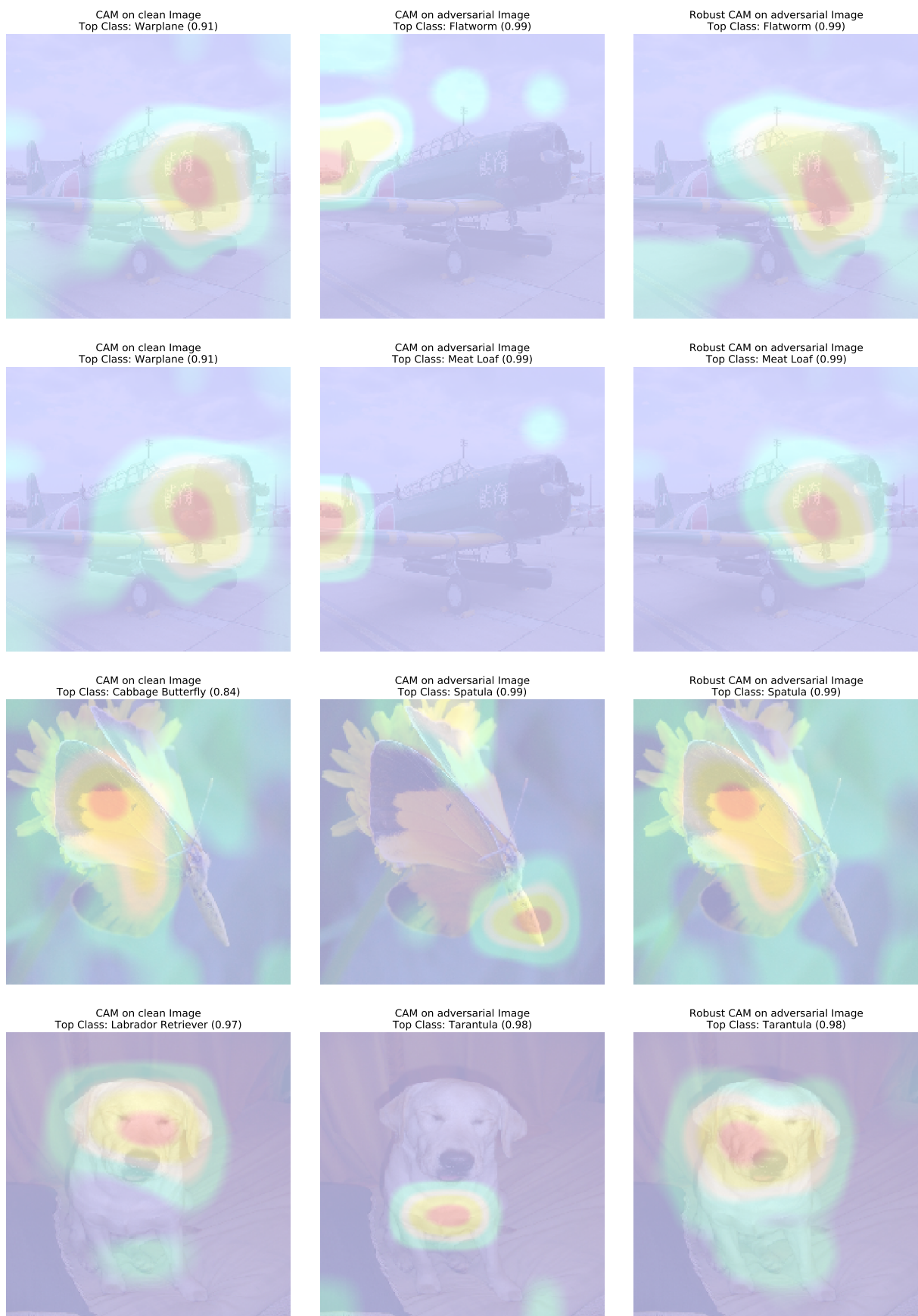


Figure 2: Comparison of Class activation maps and Robust Activation maps

Full size figures of Figure 5 (Linear search for model parameters on training data)

