
Interpreting Adversarially Trained Convolutional Neural Networks

Tianyuan Zhang¹ Zhanxing Zhu^{2,3,4}

Abstract

We attempt to interpret how adversarially trained convolutional neural networks (AT-CNNs) recognize objects. We design systematic approaches to interpret AT-CNNs in both qualitative and quantitative ways and compare them with normally trained models. Surprisingly, we find that adversarial training alleviates the texture bias of standard CNNs when trained on object recognition tasks, and helps CNNs learn a more shape-biased representation. We validate our hypothesis from two aspects. First, we compare the saliency maps of AT-CNNs and standard CNNs on clean images and images under different transformations. The comparison could visually show that the prediction of the two types of CNNs is sensitive to dramatically different types of features. Second, to achieve quantitative verification, we construct additional test datasets that destroy either textures or shapes, such as style-transferred version of clean data, saturated images and patch-shuffled ones, and then evaluate the classification accuracy of AT-CNNs and normal CNNs on these datasets. Our findings shed some light on why AT-CNNs are more robust than those normally trained ones and contribute to a better understanding of adversarial training over CNNs from an interpretation perspective.

1. Introduction

Convolutional neural networks (CNNs) have achieved great success in a variety of visual recognition tasks (Krizhevsky et al., 2012; Girshick et al., 2014; Long et al., 2015) with their stacked local connections. A crucial issue is to understand what is being learned after training over thousands or even millions of images. This involves interpreting CNNs.

¹School of EECS, Peking University, China ²School of Mathematical Sciences, Peking University, China ³Center for Data Science, Peking University ⁴Beijing Institute of Big Data Research. Correspondence to: Zhanxing Zhu <zhanxing.zhu@pku.edu.cn>.

Along this line, some recent works showed that standard CNNs trained on ImageNet make their predictions rely on the local textures rather than long-range dependencies encoded in the shape of objects (Geirhos et al., 2019; Brendel & Bethge, 2019; Ballester & de Araújo, 2016). Consequently, this texture bias prevents the trained CNNs from generalizing well on those images with distorted textures but maintained shape information. Geirhos et al. (2019) also showed that using a combination of Stylized-ImageNet and ImageNet can alleviate the texture bias of standard CNNs. It naturally raises an intriguing question:

Are there any other trained CNNs are more biased towards shapes?

Recently, normally trained neural networks were found to be easily fooled by maliciously perturbed examples, i.e., adversarial examples (Goodfellow et al., 2014; Kurakin et al., 2016). To defense the adversarial examples, adversarial training was proposed; that is, instead of minimizing the loss function over the clean example, it minimizes almost worst-case loss over the slightly perturbed examples (Madry et al., 2018). We name these adversarially trained networks as AT-CNNs. They were extensively shown to be able to enhance the robustness, i.e., improving the classification accuracy over the adversarial examples. Then,

What is learned by adversarially trained CNNs to make it more robust?

In this work, in order to explore the answer to the above questions, we systematically design various experiments to interpret the AT-CNNs and compare them with normally trained models. We find that AT-CNNs are better at capturing long-range correlations such as shapes, and less biased towards textures than normally trained CNNs in popular object recognition datasets. This finding partially explains why AT-CNNs tends to be more robust than standard CNNs.

We validate our hypothesis from two aspects. First, we compare the saliency maps of AT-CNNs and standard CNNs on clean images and those under different transformations. The comparison could visually show that the predictions of the two CNNs are sensitive to dramatically different types of features. Second, we construct additional test datasets that destroy either textures or shapes, such as the style-transferred version of clean data, saturated images and patch-

shuffled images, then evaluate the classification accuracy of AT-CNN and normal CNNs on these datasets. These sophisticated designed experiments provide a quantitative comparison between the two CNNs and demonstrate their biases when making predictions.

To the best of our knowledge, we are the first to implement systematic investigation on interpreting the adversarially trained CNNs, *both visually and quantitatively*. Our findings shed some light on why AT-CNNs are more robust than those normally trained ones and also contribute to better understanding adversarial training over CNNs from an interpretation perspective.

The remaining of the paper is structured as follows. We introduce background knowledge on adversarial training and salience methods in Section 2. The methods for interpreting AT-CNNs are described in Section 3. Then we present the experimental results to support our findings in Section 4. The related works and discussions are presented in Section 5. Section 6 concludes the paper.

2. Preliminary

2.1. Adversarial training

This training method was first proposed by (Goodfellow et al., 2014), which is the most successful approach for building robust models so far for defending adversarial examples (Madry et al., 2018; Sinha et al., 2018; Athalye et al., 2018). It can be formulated as solving a robust optimization problem (Shaham et al., 2015)

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{S}} \ell(f(x + \delta; \theta), y) \right], \quad (1)$$

where $f(x; \theta)$ represents the neural network parameterized by weights θ ; the input-output pair (x, y) is sample from the training set \mathcal{D} ; δ denotes the adversarial perturbation and $\ell(\cdot, \cdot)$ is the chosen loss function, e.g. cross entropy loss. \mathcal{S} denotes a certain norm constraints, such as l_{∞} or l_2 .

The inner maximization is approximated by adversarial examples generated by various attack methods. Training against a projected gradient descent (PGD, Madry et al. (2018)) adversary leads to state-of-the-art white-box robustness. We use PGD based adversarial training with bounded l_{∞} and l_2 norm constraints. We also investigate FGSM (Goodfellow et al., 2014) based adversarial training.

2.2. Salience maps

Given a trained neural network, visualizing the salience maps aims at assigning a *sensitivity* value, sometimes also called ‘‘attribution’’, to show the sensitivity of the output to each pixel of an input image. Salience methods can mainly be divided into (Ancona et al., 2018) *perturbation-*

based methods (Zeiler & Fergus, 2014; Zintgraf et al., 2017) and *gradient-based* method (Erhan et al., 2009; Simonyan et al., 2013; Shrikumar et al., 2017; Sundararajan et al., 2017; Selvaraju et al., 2017; Zhou et al., 2016; Smilkov et al., 2017; Bach et al., 2015). Recently (Adebayo et al., 2018) carries out a systematic test for many of the gradient-based salience methods, and only variants of Grad and GradCAM (Selvaraju et al., 2017) pass the proposed sanity checks. We thus choose Grad and its smoothed version SmoothGrad (Smilkov et al., 2017) for visualization.

Formally, let $x \in \mathbb{R}^d$ denote the input image, a trained network is a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^K$, where K is the total number of classes. Let S_c denotes the class activation function for each class c . We seek to obtain a salience map $E \in \mathbb{R}^d$. The **Grad** explanation is the gradient of class activation with respect to the input image x ,

$$E = \frac{\partial S_c(x)}{\partial x}. \quad (2)$$

SmoothGrad (Smilkov et al., 2017) was proposed to alleviate noises in gradient explanation by averaging over the gradient of noisy copies of an input. Thus for an input x , the smoothed variant of Grad, SmoothGrad can be written as

$$E = \frac{1}{n} \sum_{i=1}^n \frac{\partial S_c(x_i)}{\partial x_i}, \quad (3)$$

where $x_i = x + g_i$, and g_i are noise vectors drawn i.i.d from a Gaussian distribution $\mathcal{N}(0, \sigma^2)$. In all our experiments, we set $n = 100$, and the noise level, $\sigma / (x_{max} - x_{min}) = 0.1$. We choose $S_c(x) = \log p_c(x)$, where $p_c(x)$ is the probability of class c assigned by a classifier to input x .

3. Methods

In this section, we elaborate our method for interpreting the adversarially trained CNNs and comparing them with normally trained ones. Three image datasets are considered, including Tiny ImageNet¹, Caltech-256 (Griffin et al., 2007) and CIFAR-10.

We first visualize the salience maps of AT-CNNs and normal CNNs to demonstrate that the two models trained with different ways are sensitive to different kinds of features. Besides this qualitative comparison, we also test the two kinds of CNNs on different transformed datasets to distinguish the difference of their preferred features.

3.1. Visualizing the salience maps

A straightforward way of investigating the difference between AT-CNNs and CNNs is to visualize which group of pixels the network outputs are most sensitive to. Salience

¹<https://tiny-imagenet.herokuapp.com/>

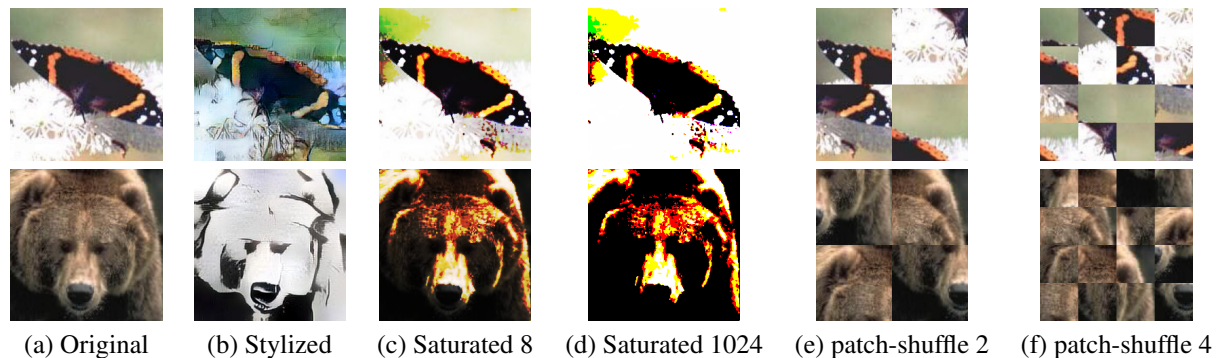


Figure 1. Visualization of three transformations. Original images are from Caltech-256. From left to right, original, stylized, saturation level as 8, 1024, 2×2 patch-shuffling, 4×4 patch-shuffling.

maps generated by Grad and its smoothed variant SmoothGrad are good candidates to show what features a model is sensitive to. We compare the saliency maps between AT-CNNs and CNNs on clean images, and images under texture preserving and shape preserving distortions. Extensive results can be seen in Section 4.1.

As pointed by Smilkov et al. (2017), sensitivity maps based on Grad method are often visually noisy, highlighting that some pixels, to a human eye, seem randomly selected. SmoothGrad in Eq. (3), on the other hand, could reduce visual noise by averaging the gradient over the Gaussian perturbed images. Thus, we mainly report the saliency maps produced by SmoothGrad, and the Grad visualization results are provided in Supplementary Materials. Note that the two visualization methods could help us draw a consistent conclusion on the difference between the two trained CNNs.

3.2. Generalization on shape/texture preserving distortions

Besides visual inspection of sensitivity maps, we propose to measure the sensitivity of AT-CNNs and CNNs to different features by evaluating the performance degradation under several distortions that either preserves shapes or textures. Intuitively, if one model relies on textures a lot, the performance would degrade severely if we destroy most of the textures while preserving other information, such as the shapes and other features. However, a perfect disentanglement of texture, shape and other feature information is impossible (Gatys et al., 2015). In this work, we mainly construct three kinds of image translations to achieve the shape or texture distortion, style-transfer, saturating and patch-shuffling operation. Some of the image samples are shown in Figure 1. We also added three Fourier-filtered test set in the supplementary files. We now describe each of these transformations and their properties.

Note that we conduct normal training or adversarial training on the original training sets, and then evaluate their general-

izability over the transformed data. During the training, we never use the transformed datasets.

Stylizing. Geirhos et al. (2019) utilized style transfer (Huang & Belongie, 2017) to generate images with conflicting shape and texture information to demonstrate the texture bias of ImageNet-trained standard CNNs. Following the same rationale, we utilize style transfer to destroy most of the textures while preserving the global shape structures in images, and build a stylized test dataset. Therefore, with similar generalization error, models capturing shapes better should also perform better on stylized test images than those biased towards textures. The style-transferred image samples are shown in Figure 1(b).

Saturation. Similar to (Ding et al., 2019), we denote the saturation of the image x by x^p , where p indicates the saturation level ranging from 0 to ∞ . When $p = 2$, the saturation operation does not change the image. When $p \geq 2$, increasing the saturation level will push the pixel values towards binarized ones, and $p = \infty$ leads to the pure binarization. Specifically, for each pixel of image x with value $v \in [0, 1]$, its corresponding saturated pixel of x^p is defined as $\text{sign}(2v - 1)|2v - 1|^{\frac{2}{p}}/2 + 1/2$. One can observe that, from Figure 1(c) and (d), increasing saturation level can gradually destroy some texture information while preserving most parts of the contour structures.

Patch-Shuffling. To destroy long-range shape information, we split images into $k \times k$ small patches and randomly rearranging the order of these patches, with $k \in \{2, 4, 8\}$. Favorably, this operation preserves most of the texture information and destroys most of the shape information. The patch-shuffled image samples are showed in Figure 1(e), (f). Note that as k increasing, more information of the original image is lost, especially for images with low resolution.

Table 1. Accuracy and robustness of all the trained models. Robustness is measured against the PGD attack with bounded l_∞ norm. Details are listed in the Supplementary Materials. Note that underfitting CNNs have similar generalization performance with some of the AT-CNNs on clean images.

	CIFAR10		TinyImageNet		Caltech 256	
	Accuracy	Robustness	Accuracy	Robustness	Accuracy	Robustness
PGD-inf: 8	86.27	44.81	54.42	14.25	66.41	31.16
PGD-inf: 4	89.17	30.85	61.85	6.87	72.22	20.10
PGD-inf: 2	91.4	39.11	67.06	1.66	76.51	7.51
PGD-inf: 1	93.40	7.53	69.42	0.18	79.11	1.70
PGD-L2: 12	85.79	34.61	53.44	14.80	65.54	31.36
PGD-L2: 8	88.01	26.88	58.21	10.03	69.75	26.19
PGD-L2: 4	90.77	13.19	64.24	3.61	74.12	14.33
FGSM: 8	84.90	34.25	66.21	0.01	70.88	20.02
FGSM: 4	88.13	25.08	63.43	0.13	73.91	15.16
Normal	94.52	0	72.02	0.01	83.32	0
Underfit	86.79	0	60.05	0.01	69.04	0

4. Experiments and analysis

Experiments setup We describe the experiment setup to evaluate the performance of AT-CNNs and standard CNNs in data distributions manipulated by above-mentioned operations. We conduct experiments on three datasets. CIFAR-10, Tiny ImageNet and Caltech-256 (Griffin et al., 2007). Note that we do not create the style-transferred and patch-shuffled test set for CIFAR-10 due to its limited resolution.

When training on CIFAR-10, we use the ResNet-18 model (He et al., 2016a;b); for data augmentation, we perform zero paddings with width as 4, horizontal flip and random crop.

Tiny ImageNet has 200 classes of objects. Each class has 500 training images, 50 validation images, and 50 test images. All images from Tiny ImageNet are of size 64×64 . We re-scale them to 224×224 and perform random horizontal flip and per-image standardization as data augmentation.

Caltech-256 (Griffin et al., 2007) consists of 257 object categories containing a total of 30607 images. Resolution of images from Caltech is much higher compared with the above two datasets. We manually split 20% of images as the test set. We perform re-scaling and random cropping following (He et al., 2016a). For both Tiny ImageNet and Caltech-256, we use ResNet-18 model as the network architecture.

Compared models, their generalization and robustness.

For all above three datasets, we train three types of AT-CNNs, they mainly differ in the way of generating adversarial examples: FGSM, PGD with bounded l_∞ norm and PGD with bounded l_2 norm, and for each attack method we train several models under different attack strengths. Details are

listed in the supplementary material. To understand whether the difference of performance degradation for AT-CNNs and standard CNNs is due to the poor generalization (Schmidt et al., 2018; Tsipras et al., 2018) of adversarial training, we also compare the AT-CNNs with an underfitting CNN (trained over clean data) with similar generalization performance as AT-CNNs. We train 11 models on each dataset. Their generalization performance on clean data, and robustness measured by PGD attack are shown in Table 1.

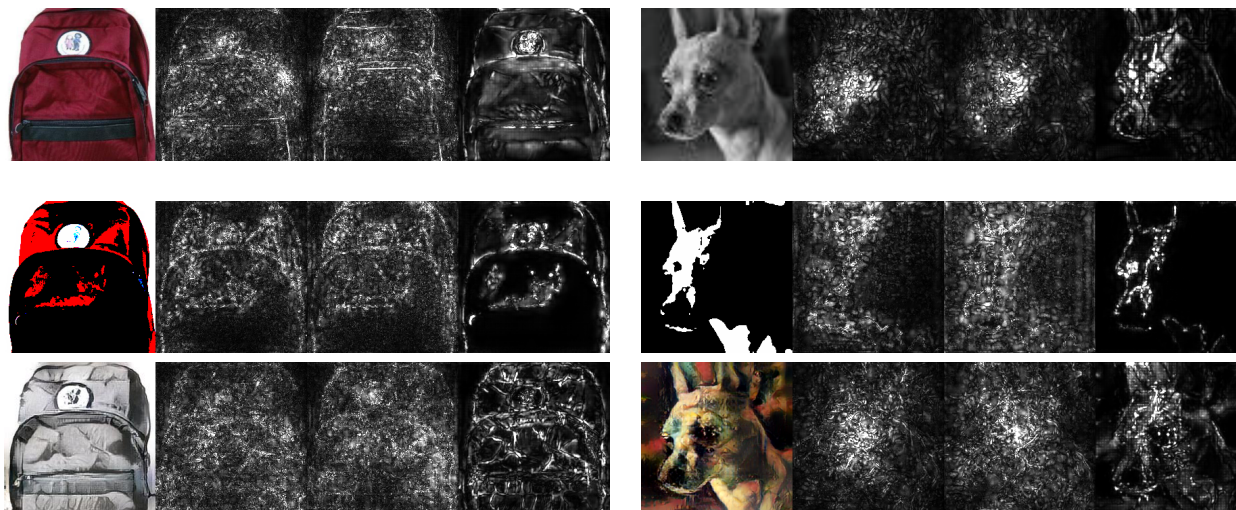
4.1. Visualization results

To investigate what features of an input image AT-CNNs and normal CNNs are most sensitive to, we generate sensitivity maps using SmoothGrad (Smilkov et al., 2017) on clean images, saturated images, and stylized images. The visualization results are presented in Figure 2.

We can easily observe that the saliency maps of AT-CNNs are much more sparse and mainly focus on contours of each object on all kinds of images, including the clean, saturated and stylized ones. Differently, sensitivity maps of standard CNNs are more noisy, and less biased towards the shapes of objects. This is consistent with the findings in (Geirhos et al., 2019).

Particularly, in the second row of Figure 2, sensitivity maps of normal CNNs of the “dog” class are still noisy even when the input saturated image are nearly binarized. On the other hand, after adversarial training, the models successfully capture the shape information of the object, providing a more interpretable prediction.

For stylized images shown in the third row of Figure 2, even with dramatically changed textures after style transfer, AT-CNNs can still be able to focus the shapes of original object,



(a) Images from Caltech-256

(b) Images from Tiny ImageNet

Figure 2. Sensitivity maps based on SmoothGrad (Smilkov et al., 2017) of three models on images under saturation, and stylizing. From top to bottom, Original, Saturation 1024 and Stylizing. For each group of images, from left to right, original image, sensitivity maps of standard CNN, underfitting CNN and PGD- l_∞ AT-CNN.

while standard CNNs totally fail.

Due to the limited space, we provide more visualization results (including the sensitivity maps generated by Grad method) in Supplementary Materials.

4.2. Generalization performance on transformed data

In this part, we mainly show generalization performance of AT-CNNs and normal CNNs on either shape or texture preserving distorted image datasets. This could help us to understand how different that the two types of models are biased in a quantitative way.

For all experimental results below, besides the top-1 accuracy, we also report an “accuracy on correctly classified images”. This accuracy is measured by first selecting the images from the clean test set that is being correctly classified, then measuring the accuracy of transformed images from these correctly classified ones.

4.2.1. STYLIZING

Following Geirhos et al. (2019), we generate stylized version of test set for Caltech-256 and Tiny ImageNet.

We report the “accuracy on correctly classified images” of all the trained models on stylized test set in Table 2. Compared with standard CNNs, though with a lower accuracy on original test images, AT-CNNs achieve higher accuracy on stylized ones with textures being dramatically changed. The comparison quantitatively shows that AT-CNNs tend to be more invariant with respect to local textures.

4.2.2. SATURATION

We use the saturation operation to manipulate the images, and show the how increasing saturation levels affects the accuracy of models trained in different ways.

In Figure 4, we visualize images with varying saturation levels. It can be easily observed that increasing saturation levels pushes images more “binnarized”, where some textures are wiped out, but produces sharper edges and preserving shape information. When saturation level is smaller than 2, i.e. clean image, it pushes all the pixels towards $1/2$ and nearly all the information is lost, and $p = 0$ leads to a totally gray image with constant pixel value.

We measure the “accuracy on correctly classified images” for all the trained models, and show them in Figure 5. We can observe that with the increasing level of saturation, more texture information is lost. Favorably, adversarially trained models exhibit a much less sensitivity to this texture loss, still obtaining a high classification accuracy. The results indicate that AT-CNNs are more robust to “saturation” or “binarizing” operations, which may demonstrate that the prediction capability of AT-CNNs relies less on texture and more on shapes. Results on CIFAR-10 tells the same story, as presented in supplementary material due to the limited space.

Additionally, in our experiments, for each adversarial training approach, either PGD or FGSM based, AT-CNNs with higher robustness towards PGD adversary are more invariant to the increasing of the saturation level and texture loss. On the other hand, adversarial training with higher robustness

Interpreting Adversarially Trained Convolutional Neural Networks



Figure 3. Visualization of images from style-transferred test set. Applying AdaIn (Huang & Belongie, 2017) style transfer distorts local textures of original images, while the global shape structure is retained. The first row are images from Caltech-256, and the second row are images from Tiny ImageNet.

Table 2. “Accuracy on correctly classified images” for different models on stylized test set. The columns named “Caltech-256” and “TinyImageNet” show the generalization of different models on the clean test set.

DATASET	CALTECH-256	STYLIZED CALTECH-256	TINYIMAGENET	STYLIZED TINYIMAGENET
STANDARD	83.32	16.83	72.02	7.25
UNDERFIT	69.04	9.75	60.35	7.16
PGD- l_∞ : 8	66.41	19.75	54.42	18.81
PGD- l_∞ : 4	72.22	21.10	61.85	20.51
PGD- l_∞ : 2	76.51	21.89	67.06	19.25
PGD- l_∞ : 1	79.11	22.07	69.42	18.31
PGD- l_2 : 12	65.24	20.14	53.44	19.33
PGD- l_2 : 8	69.75	21.62	58.21	20.42
PGD- l_2 : 4	74.12	22.53	64.24	21.05
FGSM: 8	70.88	21.23	66.21	15.07
FGSM: 4	73.91	21.99	63.43	20.22



Figure 4. Illustration of how varying saturation changes the appearance of the image. From left to right, saturation level 0.25, 0.5, 1, 2 (original image), 4, 8, 16, 64, 1024. Increasing saturation level pushes pixels towards 0 or 1, which preserves most of the shape while wiping most of the textures. Decreasing saturation level pushes all pixels to 1/2.

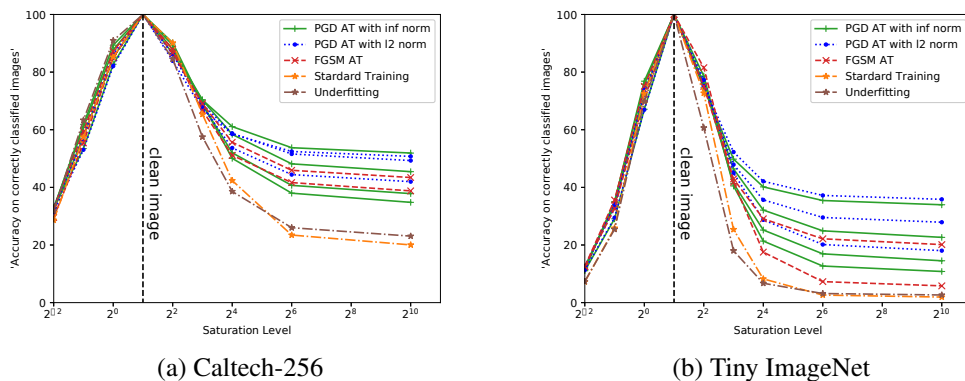


Figure 5. “Accuracy on correctly classified images” for different models on saturated Caltech-256 and Tiny ImageNet with respect to different saturation levels. Note that in the plot, there are several curves with same color and line type shown for each adversarial training method, PGD and FGSM-based, those of which with larger perturbation achieves better robustness for most of the cases. Detailed results are list in supplementary materials.

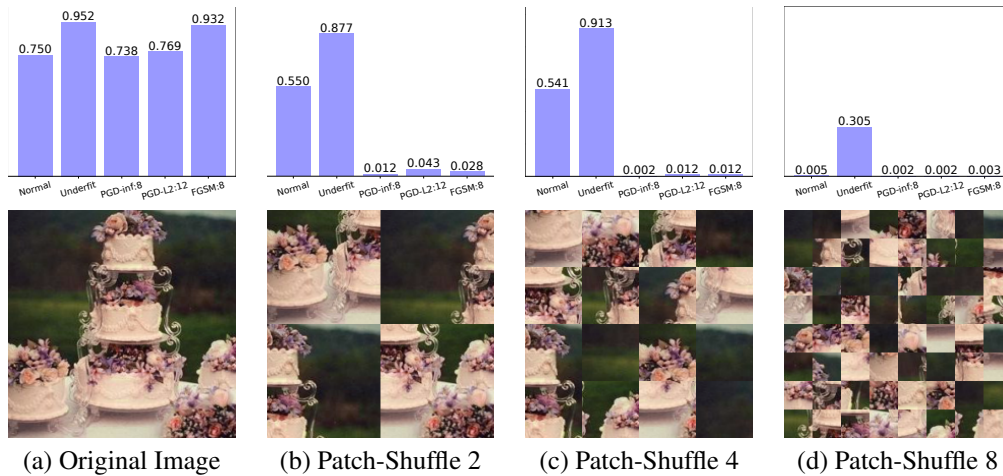


Figure 6. Visualization of patch-shuffling transformation. The first row shows probability of “cake” assigned by different models.

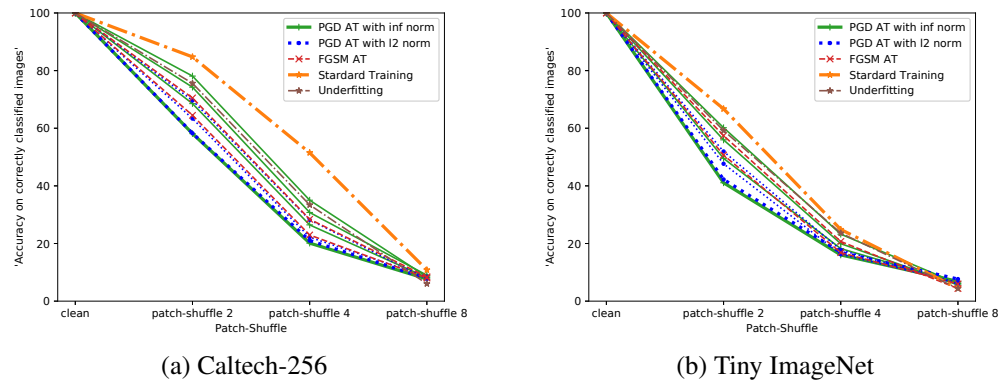


Figure 7. “Accuracy on correctly classified images” for different models on patch-shuffled Tiny ImageNet and Caltech-256 with different splitting numbers. Detailed results are listed in supplementary materials.

typically ruin the generalization over the clean dataset. Our finding also supports the claim “robustness maybe at odds with accuracy” (Tsipras et al., 2018).

When decreasing the saturation level, all models have similar degree of performance degradation, indicating that AT-CNNs are not robust to all kinds of image distortions. They tend to be more robust for fixed types of distortions. We leave the further investigation regarding this issue as future work.

4.2.3. PATCH-SHUFFLING

Stylizing and saturation operation aim at changing or removing the texture information of original images, while preserving the features of shapes and edges. In order to test the different bias of AT-CNN and standard CNN in the other way around, we shatter the shape and edge information by splitting the images into $k \times k$ patches and then randomly shuffling them. This operation could still maintains the local textures if k is not too large.

Figure 6 shows one example of patch-shuffled images under different numbers of splitting. The first row shows the probabilities assigned by different models to the ground truth class of the original image. Obviously, after random shuffling, the shapes and edge features are destroyed dramatically, the prediction probability of the adversarially trained CNNs drops significantly, while the normal CNNs still maintains a high confidence over the ground truth class. This reveals AT-CNNs are more biased towards shapes and edges than normally trained ones.

Moreover, Figure 7 depicts the “accuracy of correctly classified images” for all the models measured on “Patch-shuffled” test set with increasing number of splitting pieces. AT-CNNs, especially trained against with a stronger attack are more sensitive to “Patch-shuffling” operations in most of our experiments.

Note that under “Patch-shuffle 8” operation, all models have similar “accuracy of correctly classified images”, which is largely due to the severe information loss. Also note that this

accuracy of all models on Tiny ImageNet shown in 7(a) is much lower than that on Caltech-256 in 7(b). That is, under “Patch-shuffle 1”, normally trained CNN has an accuracy of 84.76% on Caltech-256, while only 66.73% on Tiny ImageNet. This mainly originates from the limited resolution of Tiny ImageNet, since “Patch-Shuffle” operation on low-resolution images destroys more useful features than those with higher resolution.

5. Related work and discussion

Interpreting AT-CNNs. Recently there are some relevant findings indicating that AT-CNNs learn fundamentally different feature representations than standard classifiers. Tsipras et al. (2018) showed that sensitivity maps of AT-CNNs in the input space align well with human perception. Additionally, by visualizing large- ϵ adversarial examples against AT-CNNs, it can be observed that the adversarial examples could capture salient data characteristics of a different class, which appear semantically similar to the images of the different class. Dong et al. (2017) leveraged adversarial training to produce a more interpretable representation by visualizing active neurons. Compared with Tsipras et al. (2018) and Dong et al. (2017), we have conducted a more systematical investigation for interpreting AT-CNNs. We construct three types of image transformation that can largely change the textures while preserving shape information (i.e. stylizing and saturation), or shatter the shape/edge features while keeping the local textures (i.e. patch-shuffling). Evaluating the generalization of AT-CNNs over these designed datasets provides a *quantitative* way to verify and interpret their strong shape-bias compared with normal CNNs.

Insights for defending adversarial examples. Based on our investigation over the AT-CNNs, we find that the robustness towards adversarial examples is correlated with the capability of capturing long-range features like shapes or contours. This naturally raises the question: *whether any other models that can capture more global features or with more texture invariance could lead to more robustness to adversarial examples, even without adversarial training?* This might provide us some insights on designing new network architecture or new strategies for enhancing the bias towards long-range features. Some recent works turn out partially answering this question. (Xie et al., 2018) enhanced standard CNNs with non-local blocks inspired from (Wang et al., 2018; Vaswani et al., 2017) which capture long-range dependencies in a data-dependent manner, and when combined with adversarial training, their networks achieved state-of-the-art adversarial robustness on ImageNet. (Luo et al., 2018) destroyed some of the local connection of standard CNNs by randomly select a set of neurons and remove them from the network before training, and thus forcing the CNNs to less focus on local texture features. With this

design, they achieved improved black-box robustness.

Adversarial training with other types of attacks. In this work, we mainly interpret the AT-CNNs based on norm-constrained perturbation over the original images. It is worthy of noting that the difference between normally trained and adversarially trained CNNs may highly depends on the type of adversaries. Models trained against spatially-transformed adversary (Xiao et al., 2018), denoted as ST-ST-CNNs, have similar robustness towards PGD attack with standard models, and their saliency maps are still quite different as shown in Figure 8. Also the average distance between saliency maps is close to that of standard CNN, which is much higher than that of PGD-AT-CNN. There exists a variety of generalized types of attacks, $x_{adv} = G(x; w)$ parameterized by w , such as spatially transformed (Xiao et al., 2018) and GAN-based adversarial examples (Song et al., 2018). We leave interpreting the AT-CNNs based on these generalized types of attacks as future work.



Figure 8. Sensitivity maps based on SmoothGrad (Smilkov et al., 2017) of three models. From left to right, original image, sensitivity maps of standard CNN, PGD- l_∞ AT-CNN and ST-AT-CNN.

6. Conclusion

From both qualitative and quantitative perspectives, we have implemented a systematic study on interpreting the adversarially trained convolutional neural networks. Through constructing distorted test sets either preserving shapes or local textures, we compare the sensitivity maps of AT-CNNs and normal CNNs on the clean, stylized and saturated images, which visually demonstrates that AT-CNNs are more biased towards global structures, such as shapes and edges. More importantly, we evaluate the generalization performance of the two models on the three constructed datasets, stylized, saturated and patch-shuffled ones. The results clearly indicate that AT-CNNs are less sensitive to the texture distortion and focus more on shape information, while the normally trained CNNs the other way around.

Understanding what a model has learned is an essential topic in both machine learning and computer vision. The strategies we propose can also be extended to interpret other neural networks, such as models for object detection and semantic segmentation.

Acknowledgement

This work is supported by National Natural Science Foundation of China (No.61806009), Beijing Natural Science Foundation (No.4184090), Beijing Academy of Artificial Intelligence (BAAI) and Intelligent Manufacturing Action Plan of Industrial Solid Foundation Program (No.JCKY2018204C004). We also appreciate insightful discussions with Dinghuai Zhang and Dr. Lei Wu.

References

- Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., and Kim, B. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pp. 9525–9536, 2018.
- Ancona, M., Ceolini, E., Oztireli, C., and Gross, M. Towards better understanding of gradient-based attribution methods for deep neural networks. In *6th International Conference on Learning Representations (ICLR 2018)*, 2018.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018.
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- Ballester, P. and de Araújo, R. M. On the performance of googlenet and alexnet applied to sketches. In *AAAI*, pp. 1124–1128, 2016.
- Brendel, W. and Bethge, M. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. In *International Conference on Learning Representations*, 2019.
- Ding, G. W., Lui, K. Y.-C., Jin, X., Wang, L., and Huang, R. On the sensitivity of adversarial robustness to input data distributions. In *International Conference on Learning Representations*, 2019.
- Dong, Y., Su, H., Zhu, J., and Bao, F. Towards interpretable deep neural networks by leveraging adversarial examples. *arXiv preprint arXiv:1708.05493*, 2017.
- Erhan, D., Bengio, Y., Courville, A., and Vincent, P. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.
- Gatys, L. A., Ecker, A. S., and Bethge, M. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.
- Girshick, R., Donahue, J., Darrell, T., and Malik, J. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 580–587, 2014.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Griffin, G., Holub, A., and Perona, P. Caltech-256 object category dataset. 2007.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016a.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016b.
- Huang, X. and Belongie, S. Arbitrary style transfer in real-time with adaptive instance normalization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 1510–1519. IEEE, 2017.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Kurakin, A., Goodfellow, I., and Bengio, S. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- Long, J., Shelhamer, E., and Darrell, T. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.
- Luo, T., Cai, T., Zhang, M., Chen, S., and Wang, L. Random mask: Towards robust convolutional neural networks. 2018.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., and Madry, A. Adversarially robust generalization requires more data. *arXiv preprint arXiv:1804.11285*, 2018.

- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 618–626. IEEE, 2017.
- Shaham, U., Yamada, Y., and Negahban, S. Understanding adversarial training: Increasing local stability of neural nets through robust optimization. *arXiv preprint arXiv:1511.05432*, 2015.
- Shrikumar, A., Greenside, P., and Kundaje, A. Learning important features through propagating activation differences. *arXiv preprint arXiv:1704.02685*, 2017.
- Simonyan, K., Vedaldi, A., and Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- Sinha, A., Namkoong, H., and Duchi, J. Certifiable distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Song, Y., Shu, R., Kushman, N., and Ermon, S. Constructing unrestricted adversarial examples with generative models. In *Advances in Neural Information Processing Systems*, pp. 8322–8333, 2018.
- Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. *arXiv preprint arXiv:1703.01365*, 2017.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Wang, X., Girshick, R., Gupta, A., and He, K. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7794–7803, 2018.
- Xiao, C., Zhu, J.-Y., Li, B., He, W., Liu, M., and Song, D. Spatially transformed adversarial examples. In *International Conference on Learning Representations*, 2018.
- Xie, C., Wu, Y., van der Maaten, L., Yuille, A., and He, K. Feature denoising for improving adversarial robustness. *arXiv preprint arXiv:1812.03411*, 2018.
- Zeiler, M. D. and Fergus, R. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2921–2929, 2016.
- Zintgraf, L. M., Cohen, T. S., Adel, T., and Welling, M. Visualizing deep neural network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*, 2017.