

Intriguing Frequency Interpretation of Adversarial Robustness for CNNs and ViTs

Lu Chen[†] Han Yang[‡] Hu Wang[◊] Yuxin Cao^{*} Shaofeng Li[‡] Yuan Luo[†]

[†]Shanghai Jiao Tong University, China [‡]Southeast University, China

[◊]Mohamed bin Zayed University of Artificial Intelligence, UAE

^{*}National University of Singapore, Singapore

lu.chen@sjtu.edu.cn, {yanghan, shaofengli}@seu.edu.cn,

Hu.Wang@mbzuai.ac.ae, yuxincao@comp.nus.edu.sg, luoyuan@cs.sjtu.edu.cn

Abstract—Adversarial examples have attracted significant attention over the years, yet understanding their frequency-based characteristics remains insufficient. In this paper, we investigate the intriguing properties of adversarial examples in the frequency domain for the image classification task, with the following key findings. (1) As the high-frequency components increase, the performance gap between adversarial and natural examples becomes increasingly pronounced. (2) The model performance against filtered adversarial examples initially increases to a peak and declines to its inherent robustness. (3) In Convolutional Neural Networks, mid- and high-frequency components of adversarial examples exhibit their attack capabilities, while in Transformers, low- and mid-frequency components of adversarial examples are particularly effective. These results suggest that different network architectures have different frequency preferences and that differences in frequency components between adversarial and natural examples may directly influence model robustness. Based on our findings, we further conclude with three useful proposals that serve as a valuable reference to the AI model security community.

Index Terms—adversarial robustness, adversarial perturbations, high-frequency components

I. INTRODUCTION

Despite the fact that deep neural networks (DNNs) achieve remarkable performance in many fields [1]–[3], their counterintuitive vulnerability attracts increasing attention, both for safety-critical applications [4], [5] and the black-box mechanism of DNNs [6], [7]. DNNs have been found vulnerable to adversarial examples [8]–[11], where small perturbations on the input can easily change the predictions of a well-trained DNN with high confidence.

Since then, how to alleviate the vulnerability of DNNs so as to narrow the performance gap between adversarial/natural examples is another key issue. Existing methods including defensive distillation [12] and pixel denoising [13] have shown their limitations due to follow-up attack strategies [14] or gradient masking [15]. Amongst them, adversarial training [9], [16] and its variants [17], [18] indicate their reliable robustness and outperform [19]. Moreover, as a data augmentation method, adversarial training currently seems to rely on additional data [20], [21] to further improve robustness.

Recalling that high-frequency components can be potentially linked to adversarial examples [22]–[25], however, few

explorations discuss the relationship between frequency components and the attacking capabilities of adversarial examples, *i.e.*, the performance gap between adversarial/natural examples statistically. In this paper, we revisit adversarial and natural examples from frequency perspectives. We empirically identify that the performance gap between adversarial and natural examples becomes increasingly pronounced, since higher frequency components of input samples are introduced. Specifically, as high-frequency components of adversarial examples increase, the DNN performance initially rises and subsequently decreases to the adversarial robustness of the DNN. These findings have been validated across various model architectures, including Convolutional Neural Networks (ConvNets) and Vision Transformers (ViTs). We further confirmed the observations on the CIFAR-10, CIFAR-100 [26], and Tiny ImageNet [27] datasets, as well as across widely adopted attack methods, including the FGSM [28], C&W [29], PGD [16] and AutoAttack [19] methods. Furthermore, we conducted a statistical analysis of the frequency differences between adversarial and natural examples in both standard and adversarially-trained models, which shows that the frequency discrepancies are mainly concentrated in the mid-to-high frequency components. These findings suggest that the differences in frequency components between adversarial and natural examples may directly influence model robustness.

Contributions. In this paper, we reveal intriguing properties of adversarial examples in ConvNets and ViTs from a frequency perspective. (1) We identify that the performance gap between adversarial and natural examples becomes increasingly pronounced as higher frequency components are introduced. (2) The model performance against filtered adversarial examples initially increases to achieve the highest performance, and subsequently decreases to the model robustness. (3) We observe that in ConvNets, the mid- and high-frequency components of adversarial examples reflect their attack capabilities, whereas ViTs exhibit greater sensitivity to the low- and mid-frequency components. (4) We further propose three proposals that offer valuable insights for the AI model security community.

II. FREQUENCY COMPONENTS OF IMAGES ON DNNs

A. Preliminaries

Adversarial attack. The adversarial attack aims to find a small perturbation δ within the ϵ -neighborhood of a natural example \mathbf{x} to maximize the classification loss ℓ , which misleads the model to misclassify with high confidence.

$$\delta = \arg \max_{\delta} \ell(f_{\theta}(\mathbf{x} + \delta), y), \quad s.t. \|\delta\|_p \leq \epsilon, \quad (1)$$

where f denotes a DNN model with parameters θ , and (\mathbf{x}, y) denotes a pair of an image \mathbf{x} and its ground-truth label y .

Adversarial training. Adversarial training can be considered as a min-max optimization problem:

$$\theta = \arg \min_{\theta} \mathbb{E}_{\mathbf{x}} \left[\max_{\|\delta\|_p \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right]. \quad (2)$$

The inner maximization problem aims to find the worst-case perturbations to deceive the model, while the outer minimization problem is to optimize model parameters on adversarial examples to improve model robustness.

Discrete Fourier Transform for Images. The 2D Discrete Fourier Transform \mathcal{F} converts an two dimensional image $\mathbf{x} \in \mathbb{R}^{M \times N}$ into a complex-valued frequency signal.

$$\mathcal{F}(u, v) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \mathbf{x}(m, n) e^{-j2\pi(\frac{um}{M} + \frac{vn}{N})}, \quad (3)$$

where $\mathcal{F} : \mathbb{R}^{M \times N} \mapsto \mathbb{C}^{M \times N}$ is complex-valued function in the frequency domain. Here, (m, n) represents the coordinate of an image \mathbf{x} in the spatial domain, and $\mathbf{x}(m, n)$ denotes the pixel value. (u, v) represents the coordinate of the frequency spectrum, and $\mathcal{F}(u, v)$ denotes the complex frequency value. $\mathcal{F}(u, v)$ can be represented as its amplitude $|\mathcal{F}(u, v)|$ and phase $\phi(u, v)$, i.e., $\mathcal{F}(u, v) = |\mathcal{F}(u, v)|e^{j\phi(u, v)}$. To visualize the amplitude of Fourier spectrum, we shift the low frequency components to the center of the spectrum.

Inverse Discrete Fourier Transform. Given the frequency spectrum \mathcal{F} of an image, the two dimensional image \mathbf{x} can be recovered by applying the inversion of Fourier Transform.

$$\mathbf{x}(m, n) = \frac{1}{MN} \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} \mathcal{F}(u, v) e^{j2\pi(\frac{um}{M} + \frac{vn}{N})}. \quad (4)$$

B. Methodology

Frequency components of images affecting DNNs. We aim to investigate whether the differences in frequency components between adversarial and natural examples are associated with their performance gap on models. [22], [23] have indicated that adversarial perturbations may exhibit a higher concentration in the mid- and high-frequency components. Do these differences in mid- and high-frequency components statistically contribute to the performance gap between adversarial/natural examples?

To investigate the relationship between frequency differences and the performance gap, we employ a low-pass filter with varying bandwidth \mathcal{B} to isolate and pass only the low-frequency components of adversarial and natural examples. We focus on how the model performance against adversarial

examples evolves as the bandwidth \mathcal{B} increases (i.e., the higher frequency components of the images increase), and examine the trend of the performance gap between the filtered adversarial and natural examples. To further validate the impact of frequency components on model performance, we swap the frequency components between adversarial and natural examples and test accuracy on the the frequency-swapped images. Specifically, for merged adversarial examples, we remain the frequency components of adversarial examples within the bandwidth \mathcal{B} , and swap the frequency components outside the bandwidth \mathcal{B} from natural examples. Similarly, for merged natural examples, we swap the frequency components outside the bandwidth \mathcal{B} from adversarial examples.

For generating filtered images. To validate the impact of different frequency components of images on the model performance, we generate filtered images by applying a low-pass filter with varying bandwidth \mathcal{B} . Specifically, we define a low-pass filter $\mathcal{L}_{\mathcal{B}} = \{0, 1\}^{M \times N}$ with bandwidth \mathcal{B} as the operation that only allows low-frequency components within \mathcal{B} of an image to pass through, while removing high-frequency components outside \mathcal{B} , resulting in a blurred image.

$$\mathcal{L}_{\mathcal{B}}(u, v) = \begin{cases} 1, & r < \frac{\mathcal{B}}{2} \\ 0, & r > \frac{\mathcal{B}}{2} \end{cases}, \quad (5)$$

$$\mathbf{x}^{\text{filtered}} = \mathcal{F}^{-1}(\mathcal{L}_{\mathcal{B}} \circ \mathcal{F}_{\mathbf{x}}).$$

where we allow all frequency components of the frequency spectrum with a radius $r < \frac{\mathcal{B}}{2}$ to pass (i.e., $\mathcal{L}_{\mathcal{B}}(u, v) = 1$), while setting all frequency components outside of the circle with a radius $r > \frac{\mathcal{B}}{2}$ to zero (i.e., $\mathcal{L}_{\mathcal{B}}(u, v) = 0$). Then, we apply low-pass filters with varying bandwidth \mathcal{B} to the frequency spectrum \mathcal{F} of the images \mathbf{x} , resulting in the filtered spectrum $\mathcal{L}_{\mathcal{B}} \circ \mathcal{F}_{\mathbf{x}}$. Finally, the Inverse Discrete Fourier Transform \mathcal{F}^{-1} was performed on the filtered spectrum $\mathcal{L}_{\mathcal{B}} \circ \mathcal{F}_{\mathbf{x}}$ to obtain filtered images $\mathbf{x}^{\text{filtered}}$. We set \mathcal{B}/M as the frequency scale in experiments, where $\mathcal{B}/M = 1.0$ denotes the passed frequency range is tangent to the image.

For generating adversarial examples. We generate the adversarial examples using the widely adopted methods, including the Fast Gradient Sign Method (FGSM) [28], C&W [29], Projected Gradient Descent (PGD) [16] and AutoAttack [19] methods, to construct adversarial perturbations for each image in the test set. Specifically, we evaluate the robustness of ConvNets and ViTs against filtered adversarial examples on the CIFAR-10, CIFAR-100 [26], and Tiny ImageNet [27] datasets. For ℓ_{∞} -bounded adversarial examples, we use FGSM, PGD-20 (step size $1/255$), and AutoAttack with a maximum perturbation of $\epsilon = 8/255$. Additionally, ℓ_2 -bounded adversarial examples are generated using the C&W method with parameters $c = 100$ and $\kappa = 0$ [29].

To further evaluate the robustness on the adversarially-trained models [16], we train models using the ℓ_{∞} bounded adversarial examples generated by PGD-10 (step size $2/255$).

III. EXPERIMENTS

In this section, we investigate whether the differences in frequency components between adversarial and natural examples

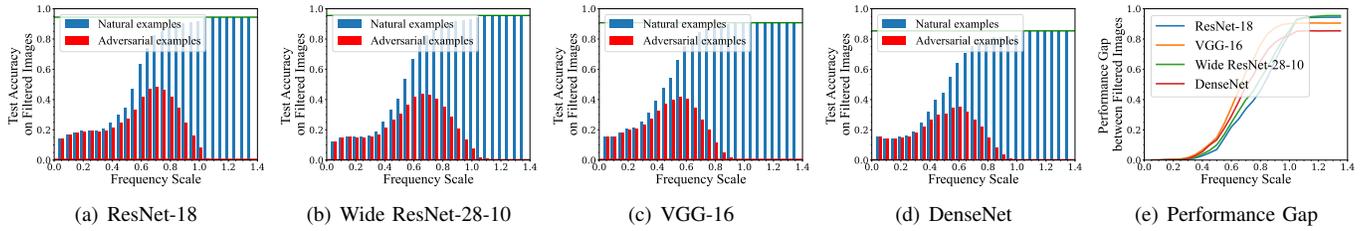


Fig. 1. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ConvNets, as evaluated on the CIFAR-10 dataset. We tested the accuracy of the filtered adversarial/natural images on (a)-(d) different ConvNets, by employing a low-pass filter with varying bandwidth on the frequency spectrum of images. (e) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases.

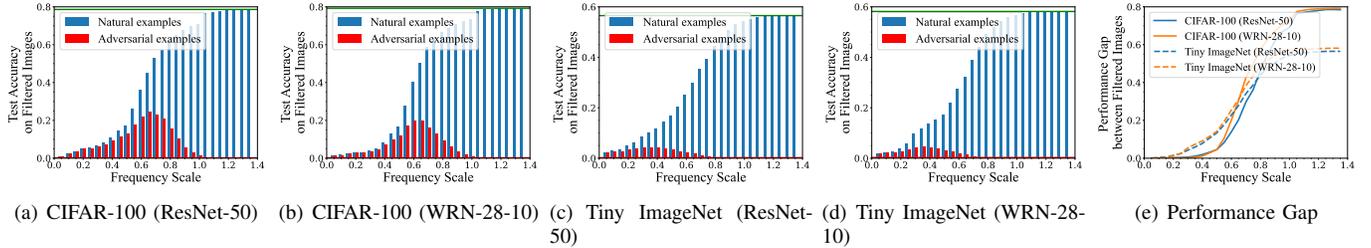


Fig. 2. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ConvNets across multiple datasets. We evaluated the accuracy of the filtered adversarial/natural images on (a)-(b) the CIFAR-100 dataset and (c)-(d) the Tiny ImageNet dataset. (e) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases across these datasets.

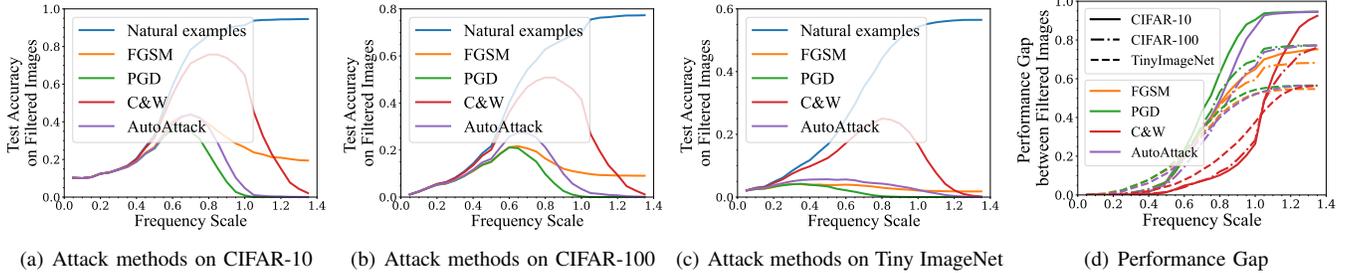


Fig. 3. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ConvNets across various attack methods. We evaluated the accuracy of the filtered adversarial examples generated using four different attack methods, the FGSM, C&W, PGD, and AutoAttack methods, on multiple datasets. (d) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases across these attack methods.

contribute to their performance gap on models.

A. Research Questions on Frequency Components of Images

We aim to investigate whether the performance of DNNs is really affected by the frequency components of adversarial and natural examples. Specifically, this section will address the following three questions:

Q1: How do the frequency components of adversarial examples affect the performance of convolutional neural networks?

Q2: What impact do the frequency components of adversarial examples have on the effectiveness of Vision Transformers?

Q3: To what extent do the frequency components of adversarial examples shape the performance of adversarially-trained models?

B. Frequency Components of images on ConvNets for Q1

We evaluate the model performance *w.r.t.* the filtered adversarial/natural examples using various ConvNets, including ResNet-18 [30], Wide ResNet-28-10 [31], VGG-16 [32], and DenseNet [33], on the CIFAR-10 dataset [26]. Fig. 1 shows

that the model performance *w.r.t.* the filtered adversarial examples generated using the PGD method, as well as natural examples, on different ConvNets. For standard ConvNets, as higher frequency components are introduced, the performance of filtered adversarial examples initially increases, peaking at an accuracy around 40%, and subsequently declines to 0.0% (red line) which reflects the DNN’s robustness against adversarial examples. In contrast, the introduction of higher frequency components in natural examples enhances classification performance, eventually reaching the clean performance of the model (green line). Fig. 1(e) shows the performance gap between filtered adversarial examples and filtered natural examples on various ConvNets. It shows that the low-frequency components of both adversarial and natural examples have a nearly identical impact on the model. However, as higher frequency components are incorporated, the performance gap between filtered adversarial and natural examples gradually increases from near zero to the gap between model generalization and robustness.

To validate the above observations, we conduct more exper-

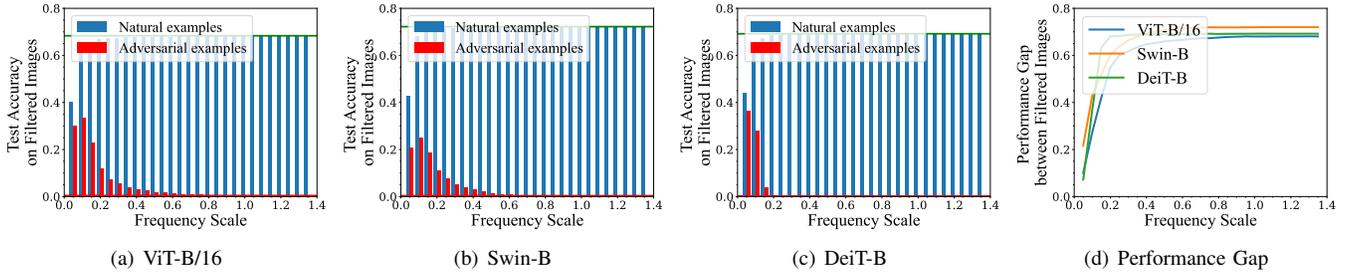


Fig. 4. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ViTs, as evaluated on the CIFAR-10 dataset. We tested the accuracy of the filtered adversarial/natural images on (a)-(c) different Transformers, by employing a low-pass filter with varying bandwidth on the frequency spectrum of images. (d) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases.

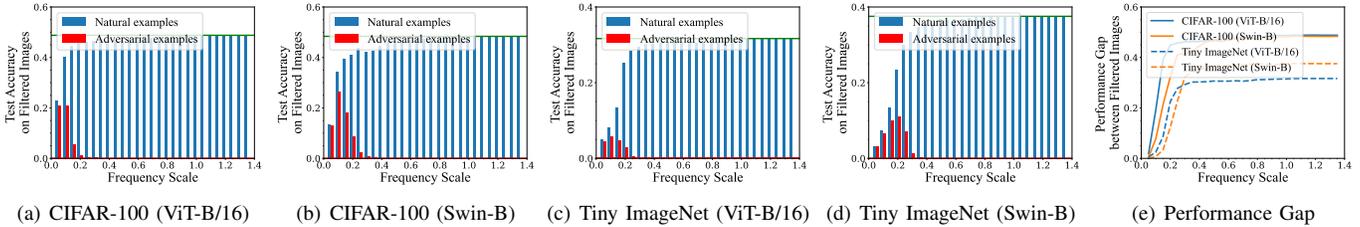


Fig. 5. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ViTs across multiple datasets. We evaluated the accuracy of the filtered adversarial/natural images on (a)-(b) the CIFAR-100 dataset and (c)-(d) the Tiny ImageNet dataset. (e) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases across these datasets.

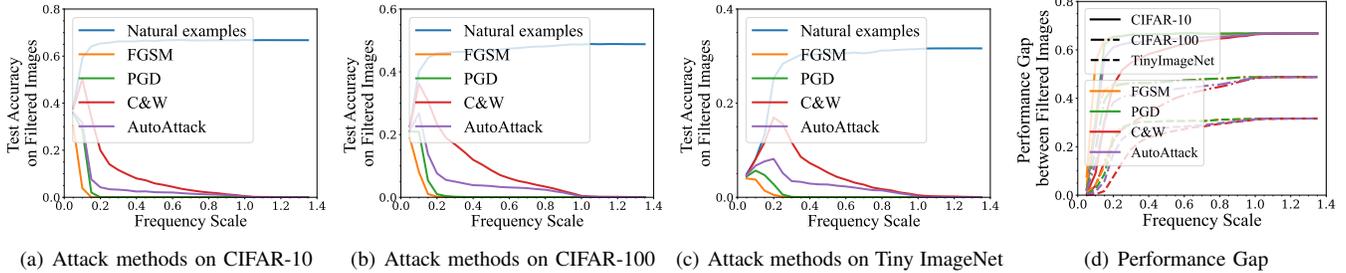


Fig. 6. The differences of frequency components contribute to the performance gap between adversarial and natural examples on ViTs across various attack methods. We evaluated the accuracy of the filtered adversarial examples generated using four different attack methods, the FGSM, C&W, PGD, and AutoAttack methods, on multiple datasets. (d) As higher frequency components are introduced, the performance gap between adversarial and natural examples increases across these attack methods.

iments across multiple datasets and various attack methods. Fig. 2 shows that the model performance *w.r.t.* the filtered adversarial and natural examples on the CIFAR-100 and Tiny ImageNet datasets, where the adversarial examples are generated using the PGD method. Similarly, the performance gap between filtered adversarial/natural examples gradually increases from zero to the gap between the model’s generalization and robustness. It is worth noting that due to variations in image size and resolution across datasets, the frequency range at which filtered adversarial examples achieve their peak performance may differ across different datasets. Fig. 3 presents that the model performance *w.r.t.* the filtered adversarial examples generated using the FGSM, C&W, PGD, and AutoAttack methods on the ResNet-50 [30]. Among these methods, FGSM and PGD are gradient-based attacks, AutoAttack is an ensemble attack that includes PGD, and C&W is an optimization-based attack. The results confirm that a similar phenomenon occurs across different attack methods.

Proposial 1: For Convolutional Neural Networks, the differences in mid- and high-frequency components between adversarial and natural examples play a critical role in their performance gap on models. These frequency components of adversarial examples exhibit their attack capabilities on models, and simply filter out these frequency components can effectively alleviate the vulnerability of models. This indicates that the researchers can focus more on the mid- and high-frequency components of the adversarial examples for detection and defence.

C. Frequency Components of images on Transformers for Q2

Fig. 4 shows that the model performance *w.r.t.* the filtered adversarial/natural examples on various Vision Transformers, including ViT-B/16 [34], Swin-B [35], and DeiT-B [36], on the CIFAR-10 dataset [26]. The adversarial examples are generated using the PGD method. Similar as the phenomenon for ConvNets, as higher frequency components are introduced, the performance of filtered adversarial examples initially increases and subsequently declines to 0.0% (red line). Compared to

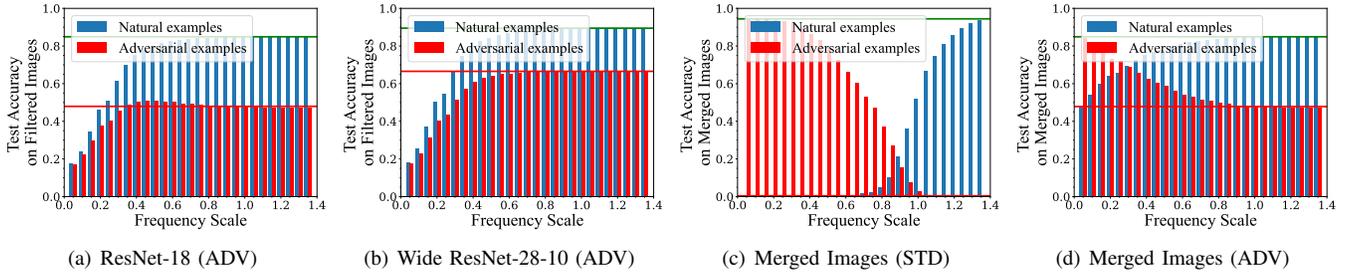


Fig. 7. The differences in frequency components statistically contribute to the performance gap between adversarial and natural examples on adversarially-trained models. We tested the accuracy of the filtered adversarial/natural images on (a)-(b) different adversarially-trained (ADV) models. We tested the accuracy of the merged adversarial/natural images on both (c) the standard (STD) model and (d) the adversarially-trained (ADV) model.

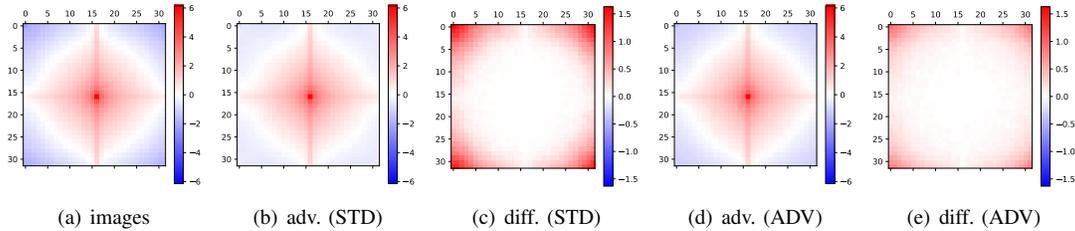


Fig. 8. The average logarithmic amplitude spectrum of (a) 1000 three-channel images and (b) their adversarial examples generated by the standard (STD) model, where the corners represent high-freq components, and the colorbars represent the logarithmic amplitude $\log(|\cdot|)$ (the redder the larger). (c) denotes the difference between (b) and (a), $\log(|adv|) - \log(|nat|) = \log(|adv|/|nat|)$, where the write color of (c) represents equivalent, and the red represents $|adv| > |nat|$. (d) and (e) are on the adversarially-trained (ADV) model, and (e) denotes the difference between (d) and (a).

ConvNets, the performance of filtered adversarial examples on Transformers declines even at low and mid frequencies, indicating that these frequency ranges contribute to the attack utility against Transformers. Fig. 4(d) further illustrates that, statistically, the low- and mid-frequency components of both adversarial and natural examples have a significant differential impact on the model performance. As higher frequency components are incorporated, the performance gap between filtered adversarial and natural examples rapidly expands from 10-20% to the gap between model generalization and robustness.

We further validate the above observations across more datasets and various attack methods on ViTs. Fig. 5 shows that the model performance *w.r.t.* the filtered adversarial and natural examples on the CIFAR-100 and Tiny ImageNet datasets. Consistent with above findings, the low- and mid-frequency components contribute to the gap between the model’s generalization and robustness across different datasets. Furthermore, Fig. 6 presents the model performance *w.r.t.* the filtered adversarial examples generated using different attacks. For Vision Transformers, adversarial examples generated by these four widely adopted attack methods predominantly exhibit their attack capabilities in the low- to mid-frequency regions.

Proposal 2: For Vision Transformers, the differences in low- and mid-frequency components between adversarial and natural examples play a critical role in their performance gap, and these frequency components of adversarial examples exhibit their attack capabilities on models. The researchers can focus more on the low- and mid-frequency components of adversarial examples to alleviate the vulnerability of Transformers, and narrow the performance gap between adversarial/natural examples.

D. Frequency Components of images on robust models for Q3

We evaluate the model performance *w.r.t.* the filtered adversarial/natural examples on adversarially-trained ResNet-18 and Wide ResNet-28-10. Fig. 7(a)-(b) show that for robust models, as higher frequency components are introduced, the model performance of filtered adversarial examples finally reaches model robustness without a rapid drop. Compared to standard models in Fig. 1(a)-(b), the low- and mid- frequency components of adversarial and natural examples initially exhibit the difference in their impact on the model performance. As higher frequency components are incorporated, the performance gap between filtered adversarial examples and natural examples increases rapidly to reach the gap between model generalization and robustness.

Comparing the frequency components of adversarial examples on standard and robust models. We further investigate the the frequency-swapped images between natural and adversarial examples on the CIFAR-10 dataset. Fig. 7(c)-(d) clearly shows the model performance *w.r.t.* the merged adversarial/natural examples on the standard and adversarially-trained ResNet-18, respectively. As higher frequency components of adversarial examples are incorporated, the classification accuracy of the merged adversarial examples gradually decreases from model generalization (green line) to model robustness (red line). Similarly, as higher frequency components of natural examples are incorporated, the classification accuracy of the merged natural examples increases from model robustness (red line) to model generalization (green line). Experiments on merged images clearly illustrate the mid- and high- frequency components of adversarial examples significantly corrupt models on standard models, while the low- and mid- frequency components of adversarial examples

corrupt models on robust models.

The statistical differences in frequency components between adversarial and natural examples. Fig. 8 visualizes the logarithmic amplitude spectrum of adversarial and natural examples. It shows that the difference in the frequency domain between two types of examples is concentrated in their high-frequency region. As Fig. 8(c) shows, compared to natural examples in Fig. 8(a), the high-frequency components of adversarial examples that generated from standard models are more pronounced in Fig. 8(b). Fig. 8(d) and 8(e) further show that high-frequency components of adversarial examples that generated from adversarially-trained models are less than those from standard models, yet still more than natural examples'. Besides, Fig. 8(e) shows that adversarial examples generated by adversarially trained models exhibit more low- and mid-frequency components.

Proposal 3: For adversarially-trained models, low- and mid-frequency components contribute to the performance gap between adversarial and natural examples. To further enhance model robustness, more attention should be directed towards improving robustness against these low- and mid-frequency components in adversarial examples.

IV. CONCLUSION

In this study, we identify intriguing properties of adversarial examples from a frequency-domain perspective. We observe the following findings. (1) The performance gap of ConvNets and ViTs between adversarial and natural examples becomes increasingly pronounced as higher frequency components are introduced. (2) The model performance against filtered adversarial examples initially increases to a peak and subsequently decreases to model robustness. (3) For ConvNets, the differences in mid- and high-frequency components contribute to their performance gap on models, whereas for ViTs the low- and mid-frequency components of adversarial examples exhibit their attack capabilities on models. Therefore, to further enhance the robustness of adversarially-trained models, more attention should be paid to the low-and mid-frequency components of adversarial examples.

REFERENCES

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*. 2019.
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *NeurIPS*. 2017.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*. 2021.
- [4] Mahmood Sharif, Sruti Bhagavatula, Lujun Bauer, and Michael K. Reiter, "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition," in *SIGSAC*. 2016.
- [5] Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Wonseok Park, Sara Rampazzi, Qi Alfred Chen, Kevin Fu, and Z. Morley Mao, "Adversarial sensor attack on lidar-based perception in autonomous driving," in *SIGSAC*. 2019.
- [6] Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George J. Pappas, "Efficient and accurate estimation of lipschitz constants for deep neural networks," in *NeurIPS*. 2019.
- [7] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry, "Adversarial examples are not bugs, they are features," in *NeurIPS*. 2019.
- [8] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus, "Intriguing properties of neural networks," in *ICLR*. 2014.
- [9] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy, "One pixel attack for fooling deep neural networks," in *ICLR*. 2015.
- [10] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard, "Universal adversarial perturbations," in *CVPR*. 2017.
- [11] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry, "Do adversarially robust imagenet models transfer better?," in *NeurIPS*. 2020.
- [12] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *SP*. 2016.
- [13] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," in *CVPR*. 2018.
- [14] Nicholas Carlini and David Wagner, "Towards evaluating the robustness of neural networks," in *SP*. 2017.
- [15] Anish Athalye, Nicholas Carlini, and David Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," in *ICML*. 2018.
- [16] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, "Towards deep learning models resistant to adversarial attacks," in *ICLR*. 2018.
- [17] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghauoui, and Michael I. Jordan, "Theoretically principled trade-off between robustness and accuracy," in *ICML*. 2019.
- [18] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu, "Improving adversarial robustness requires revisiting misclassified examples," in *ICLR*. 2020.
- [19] Francesco Croce and Matthias Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," in *ICML*. 2020.
- [20] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry, "Adversarially robust generalization requires more data," in *NeurIPS*. 2018.
- [21] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A. Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann, "Data augmentation can improve robustness," in *NeurIPS*. 2021.
- [22] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing, "High-frequency component helps explain the generalization of convolutional neural networks," in *CVPR*. 2020.
- [23] Dong Yin, Raphael Gontijo Lopes, Jonathon Shlens, Ekin D. Cubuk, and Justin Gilmer, "A fourier perspective on model robustness in computer vision," in *NeurIPS*. 2019.
- [24] Paula Harder, Franz-Josef Pfreundt, Margret Keuper, and Janis Keuper, "Spectraldefense: Detecting adversarial attacks on cnns in the fourier domain," in *IJCNN*. 2021.
- [25] Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy, "Focal frequency loss for image reconstruction and synthesis," in *ICCV*, 2021.
- [26] Alex Krizhevsky, "Learning multiple layers of features from tiny images," 2009.
- [27] Ya Le and Xuan Yang, "Tiny imagenet visual recognition challenge," 2015.
- [28] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, "Explaining and harnessing adversarial examples," *ICLR*, 2015.
- [29] Nicholas Carlini and David Wagner, "Towards evaluating the robustness of neural networks," in *SP*, 2017.
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*. 2016.
- [31] Sergey Zagoruyko, "Wide residual networks," *arXiv preprint arXiv:1605.07146*, 2016.
- [32] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [33] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger, "Densely connected convolutional networks," in *CVPR*, 2017.

- [34] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.
- [35] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021.
- [36] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou, "Training data-efficient image transformers & distillation through attention," in *ICML*, 2021.