

Supplementary Material for “Frequency Domain Model Augmentation for Adversarial Attack”

Yuyang Long¹, Qilong Zhang¹, Boheng Zeng¹, Lianli Gao¹, Xianglong Liu²,
Jian Zhang³, and Jingkuan Song^{1*}

¹ Center for Future Media, University of Electronic Science and Technology of China

² Beihang University ³ Hunan University

yuyang.long@outlook.com, {qilong.zhang, boheng.zeng}@std.uestc.edu.cn,

lianli.gao@uestc.edu.cn, xlliu@buaa.edu.cn, jianzh@hnu.edu.cn,

jingkuan.song@gmail.com

In this supplementary material, we provide: (A) theoretical proof of our proposed method, (B) discussion on hyper-parameters for our proposed S²I-FGSM, (D) discussion on spatial domain transformation analysis and visualizations for spectrum transformation images.

A Proof

Proposition 1. *Our proposed spectrum transformation can generate diverse spectrum saliency maps and thus simulate diverse substitute models.*

Proof. According to Lagrange’s mean value theorem:

$$\frac{\partial J(\mathbf{x}_1, y; \phi)}{\partial \mathbf{x}_1} = \frac{\partial J(\mathbf{x}_2, y; \phi)}{\partial \mathbf{x}_2} + \mathbf{K}, \quad (1)$$

where $\mathbf{K} = \frac{\partial^2 J(\boldsymbol{\zeta}, y; \phi)}{\partial \boldsymbol{\zeta}^2}(\mathbf{x}_1 - \mathbf{x}_2)$, $\boldsymbol{\zeta} \in [\mathbf{x}_2, \mathbf{x}_1]$.

Without spectrum transformation function $\mathcal{T}(\cdot)$, spectrum saliency map:

$$\mathbf{S}_\phi = \frac{\partial J(\mathcal{D}_\mathcal{I}(\mathcal{D}(\mathbf{x})), y; \phi)}{\partial \mathcal{D}(\mathbf{x})}, \quad (2)$$

after applying our proposed spectrum transformation function $\mathcal{T}(\cdot)$, the resulting spectrum saliency map:

$$\mathbf{S}'_\phi = \frac{\partial J(\mathcal{T}(\mathbf{x}), y; \phi)}{\partial \mathcal{D}(\mathbf{x})}, \quad (3)$$

where $\mathcal{T}(\mathbf{x}) = \mathcal{D}_\mathcal{I}((\mathcal{D}(\mathbf{x}) + \mathcal{D}(\boldsymbol{\xi})) \odot M)$

Let \mathbf{D}_1 denotes $\frac{\partial J(\mathcal{D}_\mathcal{I}(\mathcal{D}(\mathbf{x})), y; \phi)}{\partial \mathcal{D}_\mathcal{I}(\mathcal{D}(\mathbf{x}))}$ and \mathbf{D}_2 denotes $\frac{\partial \mathcal{D}_\mathcal{I}(\mathcal{D}(\mathbf{x}))}{\partial \mathcal{D}(\mathbf{x})}$, then $\mathbf{S}_\phi = \mathbf{D}_1 \mathbf{D}_2$ (according to chain rule). After applying $\mathcal{T}(\cdot)$ to \mathbf{x} , resulting spectrum saliency map \mathbf{S}'_ϕ can be expressed as:

$$\mathbf{S}'_\phi = \mathbf{D}'_1 \mathbf{D}'_2 \odot M, \quad (4)$$

*Corresponding author

where

$$\mathbf{D}'_1 = \frac{\partial J(\mathcal{D}_{\mathcal{I}}(\mathcal{D}(\mathbf{x} + \boldsymbol{\xi}) \odot \mathbf{M}), y; \phi)}{\partial \mathcal{D}_{\mathcal{I}}(\mathcal{D}(\mathbf{x} + \boldsymbol{\xi}) \odot \mathbf{M})}, \quad (5)$$

$$\mathbf{D}'_2 = \frac{\partial \mathcal{D}_{\mathcal{I}}(\mathcal{D}(\mathbf{x} + \boldsymbol{\xi}) \odot \mathbf{M})}{\partial (\mathcal{D}(\mathbf{x} + \boldsymbol{\xi}) \odot \mathbf{M})}. \quad (6)$$

Based on Eq. 1, we can formally formulate \mathbf{S}'_{ϕ} to be:

$$\begin{aligned} \mathbf{S}'_{\phi} &= (\mathbf{D}_1 + \mathbf{K}_1)(\mathbf{D}_2 + \mathbf{K}_2) \odot \mathbf{M}, \\ &= (\mathbf{S}_{\phi} + \mathbf{K}') \odot \mathbf{M}, \end{aligned} \quad (7)$$

where \mathbf{K}_1 and \mathbf{K}_2 are two specific matrices, and $\mathbf{K}' = \mathbf{D}_1\mathbf{K}_2 + \mathbf{D}_2\mathbf{K}_1 + \mathbf{K}_1\mathbf{K}_2$. Eq. 7 clearly demonstrates that our proposed transformation $\mathcal{T}(\cdot)$ is capable of simulating a different spectrum saliency map.

B On the Hyper-Parameters Settings

We first study the influence of the hyper-parameters (*i.e.*, standard deviation (std) σ of noise $\boldsymbol{\xi}$, tuning factor ρ of matrix \mathbf{M} , number N of spectrum transformations) for the proposed Spectrum Simulation Attack method.

B.1 On the Standard Deviation σ of Noise $\boldsymbol{\xi}$

In Figure 1, we report the attack success rates of S²I-FGSM for different std σ . Adversarial examples are crafted via Inc-v3 with $N = 20$ and $\rho = 0.5$. Particularly, $\sigma = 0$ means no noise is added to the input. A first glance shows that for normally trained models, the attack success rates increase gradually as σ increases and then tend to decrease when σ exceeds 16. Also when $\sigma = 16$, the defense models can achieve relatively high attack success rates. Therefore, we set $\sigma = 16$ in our paper.

B.2 On the Tuning Factor ρ of Matrix \mathbf{M}

In this section, we study the effect of tuning factor ρ for our S²I-FGSM in Figure 2. Adversarial examples are crafted via Inc-v3 with $N = 20$ and $\sigma = 16$. Particularly, $\rho = 0$ means there is no tuning on the spectrum. Similarly, as ρ increases, the degree of spectrum transformation becomes stronger and the attack success rates gradually increase and peak at $\rho = 0.5$. If we continue to increase ρ (*i.e.* $\rho > 0.5$), the attack success rates will decrease which may be attributed to the excessive spectrum transformation. To achieve better transferability, we choose $\rho = 0.5$ in our paper.

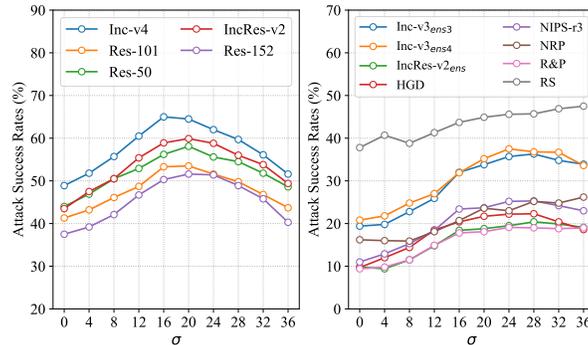


Fig. 1: The attack success rates (%) of S²I-FGSM on normally trained and defense models w.r.t. the std σ of ξ . Adversarial examples are generated via Inc-v3. **Left:** The results for fooling normally trained models. **Right:** The results for fooling defense models.

B.3 On the Number N of Spectrum Transformations.

In this section, we study the effect of number N of spectrum transformations for our S²I-FGSM in Figure 3. Adversarial examples are crafted via Inc-v3 with $\rho = 0.5$ and $\sigma = 16$. As shown in Figure 3, when $N = 1$, our method performs only one spectrum transformation and achieves the lowest transferability. As N increases, the transferability of adversarial examples is significantly enhanced at first, and turns to increase slowly after N exceeds 20. It also demonstrates that our spectrum transformation can effectively narrow the gap between the substitute model and victim model. It is worth noting that larger N implies expensive computational overhead, as we need more forward and backward propagation for gradient computation at each iteration. To balance the transferability and computational overhead, we choose $N = 20$ in our paper.

C Time Analysis of DCT/IDCT

In our experiments, we directly apply DCT/IDCT on the full image which is a time-consuming operation. Therefore, in this section we analyze the time consumption of DCT/IDCT. In Tab.1 we show the average time of an adversarial example generated by S²I-FGSM and the average time of DCT/IDCT among it. For example, let IncRes-v2 be the substitute model, S²I-FGSM takes an average of 3.78s to produce an adversarial example, of which DCT/IDCT takes up 0.58s (only accounts for 15.3% of all overheads). The experiment is conducted on RTX 3090 GPUs.

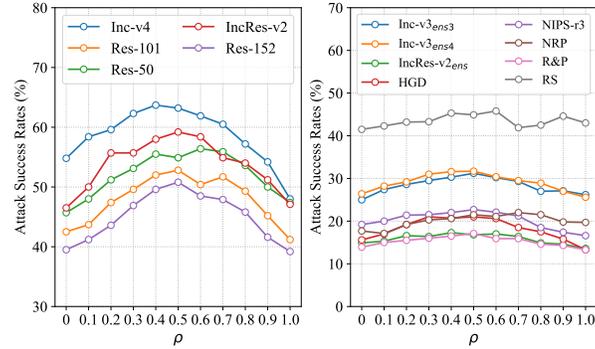


Fig. 2: The attack success rates (%) of S²I-FGSM on normally trained and defense models w.r.t. the tuning factor ρ . Adversarial examples are generated via Inc-v3. **Left:** The results for fooling normally trained models. **Right:** The results for fooling defense models.

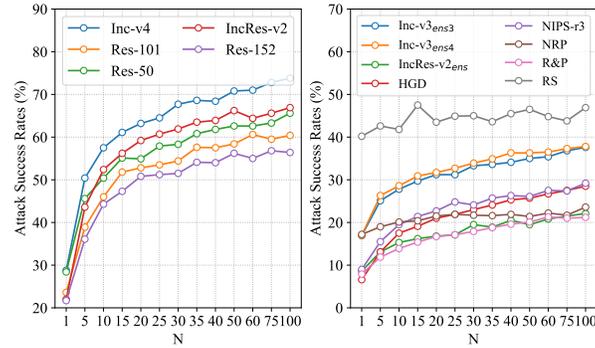


Fig. 3: The attack success rates (%) of S²I-FGSM on normally trained and defense models w.r.t. the number N of spectrum transformations. Adversarial examples are generated via Inc-v3. **Left:** The results for fooling normally trained models. **Right:** The results for fooling defense models.

Table 1: The average time (s) of generating an adversarial example on Inc-v3, Inc-v4, IncRes-v2 and Res-152, respectively. The left side of slash indicates the time of DCT/IDCT and right side indicates the time of S²I-FGSM.

| | Inc-v3 | Inc-v4 | IncRes-v2 | Res-152 |
|------|-----------|-----------|-----------|-----------|
| Time | 0.60/1.89 | 0.61/2.85 | 0.58/3.78 | 0.61/3.05 |

D Additional Results

D.1 Spatial Domain Transformation Analysis

In this section, we further validate our point that analysis on spatial domain cannot well reflect the gap between models. To support our point, we first define spatial saliency map $\hat{\mathcal{S}}_\phi$ as:

$$\hat{\mathcal{S}}_\phi = \frac{\partial J(\mathbf{x}, y; \phi)}{\partial \mathbf{x}}, \quad (8)$$

which is similar to our proposed spectrum saliency map \mathcal{S}_ϕ in Eq. 4. Then we flip the image horizontally (spatial domain transformation) and analyze their spatial saliency map and frequency saliency map. As shown in Figure 4, although spatial saliency maps between raw image and flipped image vary greatly, the changes in frequency spectrum and frequency saliency map (an indicator reflecting the characteristics of models) are small. Thus, analysis on spatial domain is unreliable and can hardly reflect the gap between models.

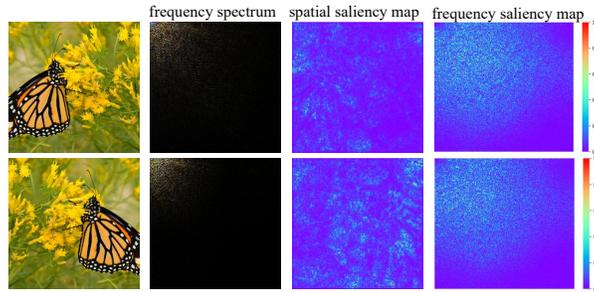


Fig. 4: Visualization for frequency spectrum, spatial saliency map, and frequency saliency map. Top row corresponds to raw image, and bottom row corresponds to spatial domain transformed image. This result demonstrates that analysis on spatial domain is unreliable.

D.2 Spectrum Transformation Images

To better understand the process of our method, we visualize the outputs of spectrum transformation. Specifically, we perform several spectrum transformations on input images and show the resulting spectrum transformation outputs in Figure 5. This figure shows that spectrum transformation just modifies colors of image and does not change its semantic information.

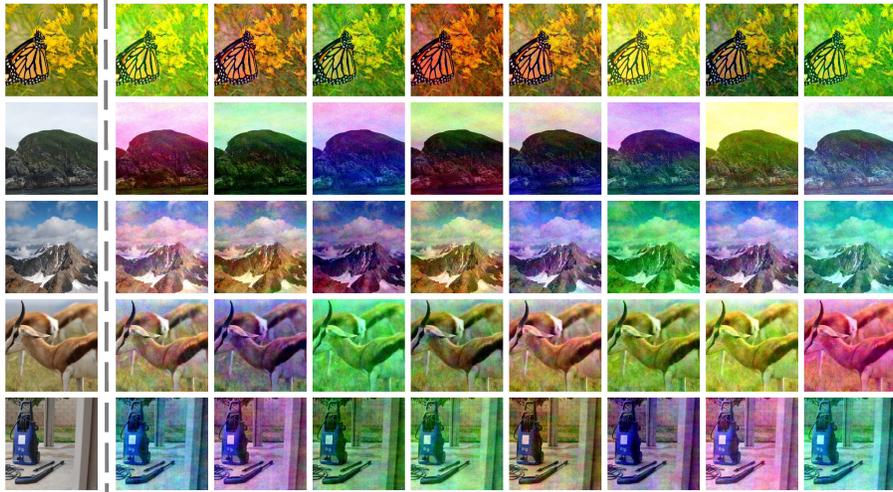


Fig. 5: Visualization for the spectrum transformation outputs (right columns) w.r.t. raw input images (left column). This result shows that spectrum transformation just modifies colors of image and does not change its semantic information.