# Unifying the Harmonic Analysis of Adversarial Attacks and Robustness

Shishira R Maiya[1], Max Ehrlich[1], Vatsal Agarwal[1], Ser-Nam Lim[2], Tom Goldstein[1], Abhinav Shrivastava[1]

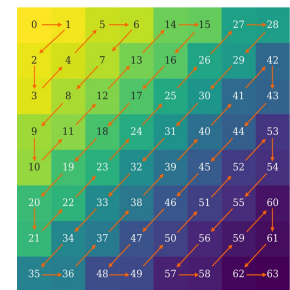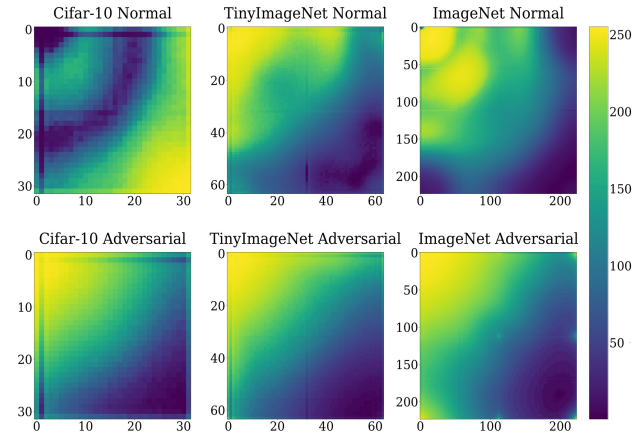[1]University of Maryland, College Park          [2]Meta

BMVC 2023

## Adversarial Examples are High frequency noise?

- Adversarial examples are imperceptible and change the output of the network when added to the input.
- The imperceptible nature makes us think they must be "High frequency noise"
- But the ineffectiveness of pre-processing methods like JPEG, deblurring and denoising as adversarial defense makes us rethink this assumption.
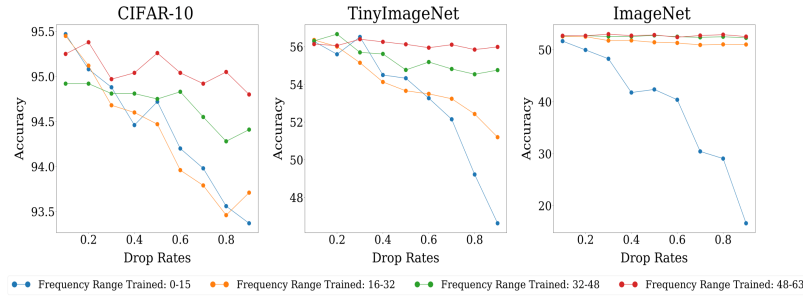
## Measuring impact of each frequency



Cifar-10 Normal, TinyImageNet Normal, ImageNet Normal

Cifar-10 Adversarial, TinyImageNet Adversarial, ImageNet Adversarial



- We plot $DCT\left(\frac{dy}{d\delta}\right)$ which measures the impact each frequency has on the resulting output predictions.
- We observe that only for CIFAR-10, higher frequencies affect the output more.
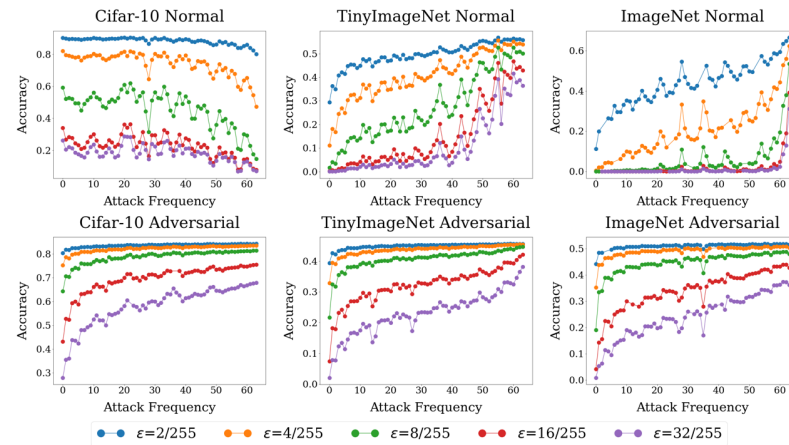- Adversarial examples are neither high nor low frequencies.

**They are dataset dependent !**

## Frequency impact during training



Frequency Range Trained: 0-15   Frequency Range Trained: 16-32   Frequency Range Trained: 32-48   Frequency Range Trained: 48-63
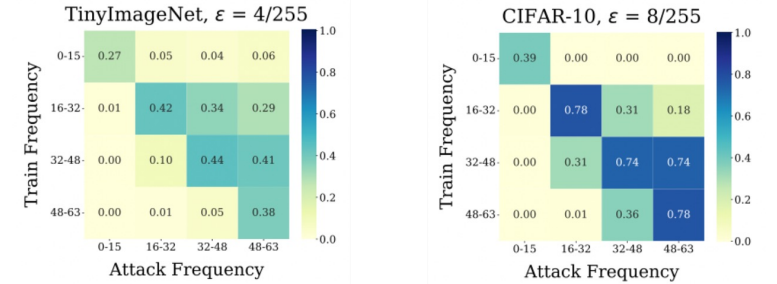
- We train models across datasets by dropping frequencies at a rate $p$ from each frequency band.
- CIFAR-10 experiences only ~2% drop when lower frequencies are dropped.
- In contrast, both ImageNet and TinyImageNet exhibit more sensitivity towards dropping of lower frequencies.

## Frequency impact on vulnerability



Cifar-10 Normal, TinyImageNet Normal, ImageNet Normal

Cifar-10 Adversarial, TinyImageNet Adversarial, ImageNet Adversarial

$\epsilon=2/255$   $\epsilon=4/255$   $\epsilon=8/255$   $\epsilon=16/255$   $\epsilon=32/255$
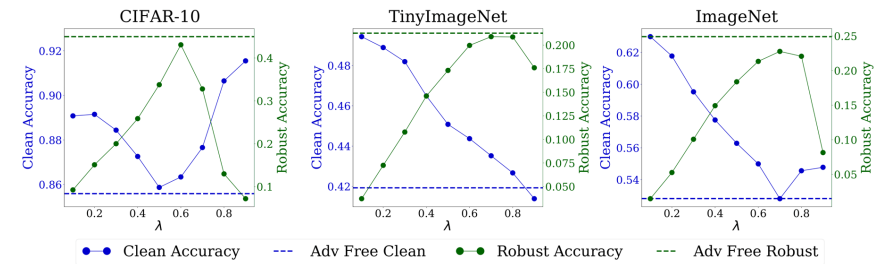
- We construct adversarial attacks by restricting them to each frequency in the DCT spectrum.
- We can observe that only for CIFAR-10 normal training, the attacks restricted on higher frequencies lead to greater reduction in accuracy.

## Adversarial Training with frequency perturbations



TinyImageNet, $\varepsilon = 4/255$

CIFAR-10, $\varepsilon = 8/255$

- Models are adversarially trained across different frequency bands and tested against other bands.
- Mid-frequency adversarial training transfers well to other bands.

## Accuracy Vs. Robustness tradeoff



CIFAR-10, TinyImageNet, ImageNet

Clean Accuracy   Adv Free Clean   Robust Accuracy   Adv Free Robust

$$\delta = \lambda \left[ \alpha \cdot \text{sgn}(\nabla_x L_{LF}) + (1 - \lambda) \cdot \left[ \alpha \cdot \text{sgn}(\nabla_x L_{HF}) \right] \right]$$

- $\lambda$ controls the amount of perturbation between low and high frequencies.
- For ImageNet and TinyImageNet, clean accuracy decreases when trained with low frequencies, while increasing robustness.
- In case of CIFAR-10, we see that there is an initial increase in robustness followed by a steep fall as higher frequencies play a significant role in robustness in this dataset.