# Impact of Spatial Frequency Based Constraints on Adversarial Robustness

Rémi Bernhard<sup>\*</sup>, Pierre-Alain Moellic<sup>\*</sup>, Martial Mermillod<sup>†</sup>, Yannick Bourrier<sup>†</sup> Romain Cohendet<sup>‡</sup>, Miguel Solinas<sup>‡</sup> Marina Reyboz<sup>‡</sup>

\* CEA-Leti, France

### <sup>†</sup> LPNC, CNRS, Université Grenoble Alpes, Université Savoie Mont Blanc, France <sup>‡</sup> Univ. Grenoble Alpes, CEA-List, France



**IJCNN 2021** 



Impact of Frequency Constraints on Adversarial Robustness



**Context:** Large-scale deployment of neural network models.

#### Attacks against integrity:

Need for a better understanding of adversarial examples, to develop efficient defenses

#### Neural computation and cognitive psychology:

 $\Rightarrow$  Importance of Low Spatial Frequency components in the human classification process

**NB:** LSF = Low Spatial Frequency, HSF = High Spatial Frequency



## **Objectives:**

- Are adversarial examples a pure HSF phenomenon?
- Link between **adversarial robustness** and **frequency properties** of information processed by models
- Investigate the robustness against adversarial perturbations offered by **spatial frequency-based constraints**

Data sets: SVHN (32x32), CIFAR10 (32x32) and "Small ImageNet" (224x224)

## Preliminaries

・ロト ・ 同ト ・ ヨト ・ ヨ

# Notations and filtering process





Figure: First row: low-pass filtering. Bottom row: high-pass filtering. For the Fourier domain masks, white  $\rightarrow$  1, black value  $\rightarrow$  0.

For LSF, low *i* value  $\Rightarrow$  strong low-pass filtering. For HSF, high *i* value  $\Rightarrow$  strong high-pass filtering.

IJCNN 2021 5 / 22

# Notations and filtering process





We filter the data set and train models on these filtered data:

X<sub>i</sub><sup>low,high</sup> → filtered data with low or high-pass filtering at intensity i (LSF/HSF task)
M<sub>i</sub><sup>low,high</sup> → Model trained on X<sub>i</sub><sup>low,high</sup>.

## Frequency properties of data and models

IJCNN 2021 7 / 22

## Frequency properties of data and models Accuracy on filtered data sets





Figure: CIFAR10 and SVHN. Accuracy of a regular model on low-pass and high-pass filtered data set.

Informative features learned (regular model):

- $\bullet~\text{SVHN} \rightarrow \text{Focused}$  on LSF task
- $\bullet~\text{CIFAR10} \rightarrow \text{Spread}$  between LSF & HSF tasks

Rémi Bernhard

IJCNN 2021 8 / 22

## Frequency properties of data and models Fourier spectrum of data





Figure: Magnitude of the Fourier spectrum for clean images. Center: LSF, border: HSF.

Relevant with frequency properties of the data sets:

- SVHN: Narrow spectrum (towards LSF)
- CIFAR10, Small ImageNet: Spread spectrum



# Frequency properties of data and models

models trained on filtered data sets



CIFAR10, SVHN: Test set accuracy of models trained on filtered data sets.

- $\bullet~\text{CIFAR10} \rightarrow$  useful information are distributed along the spectrum
- $\bullet~\text{SVHN} \rightarrow$  predominantly concentrated in the LSF.

Rémi Bernhard

# Sensitivity to HSF noise



Error rate of the model on a set of examples perturbed with noise located only in specific spatial frequencies:



Figure: CIFAR10 and SVHN. High values  $\rightarrow$  high sensitivity. Low values  $\rightarrow$  low sensitivity.

Rémi Bernhard

Impact of Frequency Constraints on Adversarial Robustness

IJCNN 2021 11 / 22

## Transferability analysis

Rémi Bernhard

Impact of Frequency Constraints on Adversarial Robustness

IJCNN 2021 12 / 22

## Transferability analysis Results





Figure: SVHN (left), CIFAR10 (middle), Small ImageNet (right). Transferability analysis.

**1)** Blue curves: Two way transferability  $M \leftarrow M_i^{low}$ 

 $\rightarrow$  The regular classification task and the LSF task share predominantly robust useful features.

## Transferability analysis Results





Figure: SVHN (left), CIFAR10 (middle), Small ImageNet (right). Transferability analysis.

- 2) Dissimilarity between the dotted and solid curves:
  - impact of non-robust features exploiting HSF

# Transferability analysis Results





Figure: SVHN (left), CIFAR10 (middle), Small ImageNet (right). Transferability analysis.

- 2) Dissimilarity between the dotted and solid curves:
  - as the high-pass filtering becomes more restrictive, the transferability  $M_i^{high} \rightarrow M$  decreases.

Rémi Bernhard

## Adversarial robustness of frequency constrained models

Adversarial robustness of frequency constrained models Objective

#### Goal:

Enforce the model to rely on useful features of the  $\ensuremath{\mathsf{LSF}}/\ensuremath{\mathsf{HSF}}$  task

### Frequency-constrained loss functions:

$$L_{i,j}^{freq}(\theta, x, y) = L^{E}(\theta, x, y) + \lambda_{1} \left\| f(x) - f(x_{i}^{low}) \right\|_{2}^{2} + \lambda_{2} \left\| f(x) - f(x_{j}^{high}) \right\|_{2}^{2}$$
$$L_{i}^{low}(\theta, x, y) = L^{E}(\theta, x, y) + \lambda_{1} \left\| f(x) - f(x_{i}^{low}) \right\|_{2}^{2}$$

#### Attack:

 $I_{\infty}$ -PGD with all sanity checks for gradient masking (false sense of security)

Rémi Bernhard



## L<sup>low</sup> results

**SVHN**: Up to 41% accuracy against PGD adversarial examples  $(L_6^{low})$  **CIFAR10**:

- No observed robustness
- 11% accuracy against PGD adversarial examples when considering  $M^{low}$  models

 $\rightarrow L^{low}$  brings robustness if:

*i)* a model relies predominantly on useful features of the LSF task (shared robust features) *ii*) it shows no sensitivity to HSF noise



L<sup>freq</sup> results

- **CIFAR10**: Up to 12% accuracy against PGD adversarial examples  $(L_{5,3}^{freq})$
- Small Imagenet: Up to 36% accuracy against PGD adversarial examples  $(L_{40,20}^{freq})$

ightarrow L<sup>freq</sup> can bring robustness in the case of information spread over the frequency spectrum

Adversarial robustness of frequency constrained models Combination with Adversarial Training

Adversarial training (reminder):

$$\delta = \underset{\|\delta\|_{\infty} \leq \epsilon}{\arg \max} \quad L^{\mathcal{E}}(\theta, x + \delta, y)$$

Resulting frequency constrained loss:

$$L_{i,j}^{AT,freq}(\theta, x, y) = L_{i,j}^{freq}(\theta, x + \delta, y)$$

Results:

SVHN: + 12% accuracy on adversarial examples  $(L_{10,4}^{AT,freq})$ CIFAR10: + 5% accuracy on adversarial examples  $(L_6^{AT,high})$ (compared with Adversarial Training with same clean accuracy)

 $\rightarrow$  Existing defense schemes can benefit from spatial frequency considerations



# Conclusion

Rémi Bernhard

Impact of Frequency Constraints on Adversarial Robustness

IJCNN 2021 21 / 22

э

・ロト ・ 同ト ・ ヨト ・ ヨ



#### Contributions:

- Adversarial examples exploit features of the whole frequency spectrum
- Models relying predominantly on useful features for the LSF task, and with a non-sensitivity to HF noise show robustness when constrained to rely on useful information for the LSF task
- when developing defenses, it is crucial to take into account the intrinsic frequency properties of data

#### **Perspectives:**

Investigate the relation between frequency properties of Adversarial Training and frequency-based constraints