# All You Need is RAW:
# Defending Against Adversarial Attacks with Camera Image Pipelines
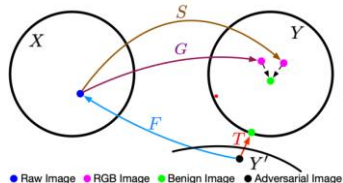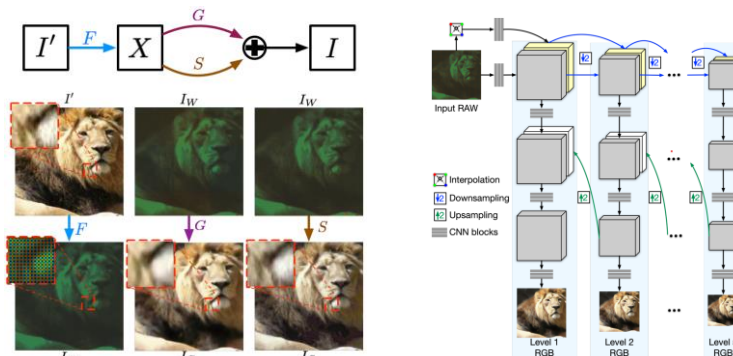
Yuxuan Zhang,  Bo Dong, Felix Heide

## Motivation

Existing neural networks for computer vision tasks are vulnerable to adversarial attacks: adding imperceptible perturbations to the input images can fool these methods to make a false prediction on an image that was correctly predicted without the perturbation. Various defense methods have proposed image-to-image mapping methods, either including these perturbations in the training process or removing them in a preprocessing denoising step. In doing so, existing methods often ignore that the natural RGB images in today's datasets are not captured but, in fact, recovered from RAW color filter array captures that are subject to various degradations in the capture. In this work, we exploit this RAW data distribution as an empirical prior for adversarial defense. Specifically, we proposed a model-agnostic adversarial defensive method, which maps the input RGB images to Bayer RAW space and back to output RGB using a learned camera image signal processing (ISP) pipeline to eliminate potential adversarial patterns. The proposed method acts as an off-the-shelf preprocessing module and, unlike model-specific adversarial training methods, does not require adversarial images to train. As a result, the method generalizes to unseen tasks without additional retraining. Experiments on large-scale datasets (e.g., ImageNet, COCO) for different vision tasks (e.g., classification, semantic segmentation, object detection) validate that the method significantly outperforms existing methods across task domains

## Model Details



Departing from existing methods, we learn a mapping from Y' to Y via an intermediate RAW distribution, X, which incorporates these RAW statistics of natural images, such as sensor photon counts, multispectral color filter array distributions and optical aberrations. To this end, the approach leverages three specially designed operators: $F: Y' \rightarrow X$, $G: X \rightarrow Y$, and $S: X \rightarrow Y$. Specifically, the F operator is a learned model, which maps an adversarial sample from its adversarial distribution to its corresponding RAW sample in the natural image distribution of RAW images. Operator G is another learned network that performs an ISP reconstruction task, i.e., it converts a RAW image to an RGB image. In theory, our goal can be achieved with these two operators by concatenating both $G(F(\cdot)): Y' \rightarrow X \rightarrow Y$. However, as these two operators are differentiable models, the potential adversary may still be able to attack the model if, under stronger attack assumptions, he has full access to the weight of preprocessing modules. To address this issue, we add the operator S, a conventional ISP, to our approach, which is implemented as a sequence of cascaded software-based sub-modules. In contrast to operator F, operator S is non-differentiable. For defending against an attack, the proposed approach first uses the F operator to map an input adversarial image, I', to its intermediate RAW measurements, IW . Then, IW is processed separately by the G and S operators to convert it to two images in the natural RGB distribution, denoted as as IG and IS, respectively. Finally, our method outputs a benign image, I, in the natural RGB distribution by combining I_G and I_S in a weighted-sum manner.

## Method Overview



● Raw Image  ● RGB Image  ● Benign Image  ● Adversarial Image

Existing defense approaches learn an RGB-to-RGB projection from an adversarial distribution (Y) to its natural RGB distribution (Y'): $T: Y' \rightarrow Y$. In contrast, our approach learns a mapping via the intermediate natural RAW distribution (X), which is achieved by utilizing three specially designed operators : $F: Y' \rightarrow X$, $G: X \rightarrow Y$, and $S: X \rightarrow Y$.

## Quantitative Result



Table 1: **Quantitative Comparisons on ImageNet** We evaluate Top-1 Accuracy on ImageNet and compare the proposed method to existing input-transformation methods. The best Top-1 accuracies are marked in bold. Our defense method offers the best performance in all settings, except for the DeepFool attack.



Table 2: **Quantitative Comparison to SOTA Input-Transformation Defense Methods on the COCO dataset.** We evaluate all methods on mean IoU (mIoU) and mark the best mIoU in bold. Our defense method offers the best performance in all settings.



Table 3: **Quantitative Comparison to SOTA Input-Transformation Defenses on the Pascal VOC dataset.** We evaluate all compared methods for mean average precision (mAP) on Pascal VOC dataset. The best mAP are marked in bold. Our defense method offers the best performance in all settings.

## Qualitative Result



**Qualitative Result:** Comparison of outputs of the proposed method along with both G and S operators and state-of-the-art defense methods on the ImageNet dataset.

## Summary

- We propose, to the best of our knowledge, the first adversarial defense method that exploits the natural distribution of RAW domain images.

- The proposed method avoids the tedious generation of adversarial training images and can be used as an off-the-shelf preprocessing module for diverse tasks.

- We validate that the method achieves state-of-the-art defense accuracy for input transformation defenses, outperforming existing approaches