

All You Need is RAW: Defending Against Adversarial Attacks with Camera Image Pipelines

Yuxuan Zhang Bo Dong Felix Heide

Princeton University

Abstract. Existing neural networks for computer vision tasks are vulnerable to adversarial attacks: adding imperceptible perturbations to the input images can fool these models 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 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 re-training. 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.

Keywords: Adversarial Defense, Low-level Imaging, Neural Image Pipeline

1 Introduction

The most successful methods for a broad range of tasks in computer vision rely on deep neural networks [11, 30, 31, 35, 91] (DNNs), including classification, detection, segmentation, scene understanding, scene reconstruction and generative tasks. Although we rely on the predictions of DNNs in safety-critical applications in robotics, self-driving vehicles, medical diagnostics, and video security, existing networks have been shown to be vulnerable to adversarial attacks [74]: small perturbations to images that are imperceptible to the human vision system to images can deceive DNNs to make incorrect predictions [52, 56, 63, 73, 78]. As such, defending against adversarial attacks [5, 51, 52, 60, 86] can help resolve

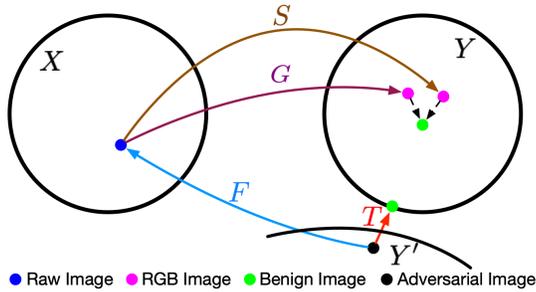


Fig. 1: Existing defense approaches learn an RGB-to-RGB projection from an adversarial distribution (Y') to its natural RGB distribution (Y): $T : Y' \rightarrow Y$. In contrast, our approach learns a mapping via the intermediate natural RAW distribution (X), which is achieved by utilizing three specially designed operators: $F : Y' \rightarrow X$, $G : X \rightarrow Y$, and $S : X \rightarrow Y$.

failure cases in safety-critical applications and provide insights into the generalization capabilities of training procedures and network architectures.

Existing defense methods fall into two approaches: They either introduce adversarial examples to the training dataset, resulting in new model weights, or they transform the inputs, aiming to remove the adversarial pattern, before feeding them into the unmodified target models. Specifically, the first line of defense methods generates adversarial examples by iteratively training a target model while finding and adding remaining adversarial images as training samples in each iteration [82, 89] [24, 80, 82]. Although the set of successful adversarial examples shrinks over time, iteratively generating them is extremely costly in training time, and different adversarial images must be included for defending against different attack algorithms. Moreover, the adversarial examples cannot be stored once in a training set as they are model-specific and domain-specific, meaning they must be re-generated when used for different models or on other domains.

Defense methods that transform the input image aim to overcome the limitations of adversarial training approaches. Considering adversarial perturbations as noise, these methods “denoise” the inputs before feeding them into unmodified target models. The preprocessing module can either employ image-to-image models such as auto-encoders or generative adversarial methods [38, 50, 67], or they rely on conventional image-processing operations [16, 19, 27, 46]. Compared to adversarial training methods, these methods are model-agnostic and require no adversarial images for training.

All methods in this approach have in *common that they rely on RGB image data as input and output*. That is, they aim to recover the distribution of natural RGB images and project the adversarial image input to the closest match in this distribution, using a direct image-to-image mapping network. As such, existing methods often ignore the fact that images in natural image datasets are the result of several processing steps applied to the raw captured images. In particular, training image datasets are produced by interpolating sub-sampled, color filtered (*e.g.*, using Bayer filter) raw data, followed by a rich low-level processing pipeline, including readout and photon noise denoising. As a result, the raw per-pixel photon counts are heavily subsampled, degraded and processed in an RGB image. We rely on the *RAW data distribution, before becoming RGB*

images, as a prior in the proposed adversarial defense method, which is empirically described in large datasets of RAW camera captures. Specifically, instead of directly learning a mapping between adversarially perturbed inputs RGBs and “clean” output RGBs, we learn a mapping via the intermediate RAW color filter array domain. In this mapping, we rely on learned ISP pipelines as low-level camera image processing blocks to map from RAW to RGB. The resulting method is entirely model-agnostic, requires no adversarial examples to train, and acts as an off-the-shelf preprocessing module that can be transferred to any task on any domain. We validate our method on large-scale datasets (ImageNet, COCO) for different vision tasks (classification, semantic segmentation, object detection), and also perform extensive ablation studies to assess the robustness of the proposed method to various attack methods, model architecture choices, and hyper-parameters choices.

Specifically, we make the following contributions:

- 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 provide a detailed analysis of how the natural RAW image distribution helps defend against adversarial attacks, and we validate that the method achieves *state-of-the-art* defense accuracy for input transformation defenses, outperforming existing approaches.

We will provide all code, models, and instructions needed to reproduce the results presented in this work.

2 Related Work

2.1 Camera Image Signal Processing (ISP) Pipeline

A camera image signal processing (ISP) pipeline converts RAW measurements from a digital camera sensor to high-quality images suitable for human viewing or downstream analytic tasks. To this end, a typical ISP pipeline encompasses a sequence of modules [39] each addressing a portion of this image reconstruction problem. In a hardware ISP, these modules are proprietary compute units, and their behavior is unknown to the user. More importantly, the modules are not differentiable [55, 76]. Two lines of works leveraged deep-learning-based approaches to cope with the significant drawback.

One line of the works directly replaced the hardware ISP with a deep-learning-based model to target different application scenarios, such as low-light enhancement [9, 10], super-resolution [87, 88, 92], smartphone camera enhancement [15, 34, 68], and ISP replacement [45]. Nevertheless, the deep-learning-based models used by these works contain a massive number of parameters and are computationally expensive. Thus, their application is limited to off-line tasks.

In contrast, another thread of works focused on searching for the best hardware ISP hyperparameters for different downstream tasks, by leveraging deep-learning-based approaches. Specifically, Tseng *et al.* [76] proposed differentiable proxy functions to model arbitrary ISP pipelines and leveraged them to find the best hardware ISP hyperparameters for different downstream tasks. Yu *et al.* [90] proposed ReconfigISP, which uses different proxy functions for each module of a hardware ISP instead of the whole ISP pipeline. Mosleh *et al.* [55] proposed a hardware-in-the-loop method to optimize hyperparameters of a hardware ISP directly.

2.2 Adversarial Attack Methods

Adversarial attacks have drawn significant attention from the deep-learning community. Based on the access level to target networks, adversarial attacks can be broadly divided into white-box attacks and black-box attacks.

Among the white-box attack, one important stream is gradient-based attacks [25, 42, 52]. These approaches generate adversarial samples based on the gradient of the loss function with respect to input images. Another flavor of attacks is based on solving optimization problems to generate adversarial samples [7, 73]. In the black-box setting, only benign images and their class labels are given, meaning attackers can only query the target model. Black-box attacks mainly leverage the free query and adversarial transferability to train substitute models [32, 59, 61, 71] or directly estimate the target model gradients [13, 14, 79] to generate adversarial examples. To avoid the transferability assumption and the overhead of gathering data to train a substitute model, several works proposed local-search-based black-box attacks to generate adversarial samples directly in the input domain [6, 44, 57].

In the physical world, adversarial samples are captured by cameras as inputs to target networks, involving camera hardware ISPs and optical systems. A variety of strategies have been developed to guard the effectiveness of the adversarial patterns in the wild [3, 18, 22, 36]. These methods typically assume that the camera acquisition and subsequent hardware processing do not alter the adversarial patterns. However, Phan *et al.* [62] have recently realized attacks of individual camera types by exploiting slight differences in their hardware ISPs and optical systems.

2.3 Defense Methods

In response to adversarial attack methods, there have been significant efforts in constructing defenses to counter those attacks. These include adversarial training [52], input transformation [4, 20], defensive distillation [60], dynamic models [83], loss modifications [58], model ensemble [69] and robust architecture [28]. Note that with the ongoing intense arms race between attacks and defenses, no defense methods are immunized to all existing attacks [1]. We next analyze the two representative categories of defense methods:

Adversarial Training (AT): The idea of AT is the following: in each training loop, it augments training data with adversarial examples generated by different attacks. AT is known to “overfit” to the attacks “seen” during training

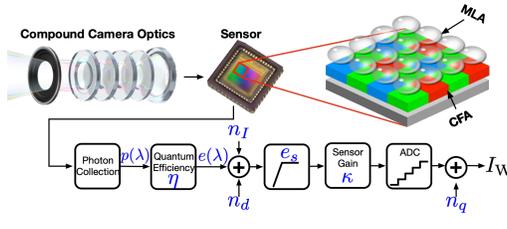


Fig. 2: Overview of the RAW imaging pipeline model. The scene light field is captured by compound camera optics, and then it is gathered by an MLA layer and fed through a CFA layer. The color-filtered photons are converted into electrons based on quantum efficiency before adding dark current and noise. Next, the converted electrons are clipped based on the maximum well capacity, e_s , and scaled by a sensor gain factor κ . Finally, an ADC converts the analog signal into a digital readout with quantization noise n_q , I_W .

and has been demonstrated to be vastly effective in defending those attacks. However, AT does not generalize well on “unseen” attacks [72]. Furthermore, iteratively generating adversarial images is time-consuming, taking 3-30 times longer than standard training before the model converges [70]. Multiple methods have been proposed to reduce the training time, making AT on large datasets (*e.g.*, ImageNet) possible [24, 80, 82, 93]. Even so, for each specific model, it still requires an extra adversarial training process and suffers from cross-domain attacks. Besides the target model, it is also worth noting that adversarial examples can be used to train the input preprocessing models. [46].

Input Transformation (IT): IT, as an image pre-processing approach, aims to remove adversarial patterns to counter attacks. A considerable number of IT methods have been proposed such as JPEG compression [16, 20, 49], randomization [84], image quilting [27], pixel deflection [64], and deep-learning-based approaches [38, 50, 67]. These IT methods can seamlessly work with different downstream models and tasks. More importantly, the IT methods can be easily combined with other model-specific defense methods to offer a stronger defense.

Our work falls into the IT category. Unlike the existing IT methods to focus on the preprocessing in the RGB distribution, the proposed approach leverages intermediate natural RAW distribution to remove adversarial patterns, which is the first work to exploit RAW distribution in the adversarial defense domain.

3 Sensor Image Formation

In this section, we review how a RAW image is formed. In short, when light from the scene enters a camera aperture, it first passes through compound camera optics. Following that is the aperture and shutter, which can be adjusted to define f-number and exposure time. Then the light falls on image sensors (*e.g.*, CCD and CMOS), where the photons are color-filtered and converted into electrons. Finally, the electrons are converted to digital values, comprising a RAW image. We refer the reader to Karaimer and Brown [40] for a detailed review.

Compound Camera Optics: A compound lens consisting of a sequence of optics is designed to correct optical aberrations. When a scene radiance, I_{SCENE} (in the form of a light field) enters a compound lens, the radiance is modulated by the complex optical pipelines and generates the image I_O , that appears on

an image sensor surface. Compound optics can be modeled by spatially-varying point spread functions (PSFs) [75].

Color Image Sensor Model: A conventional color image sensor has three layers. On the top is a micro-lens array (MLA) layer; the bottom is a matrix of small potential wells; a color filter array (CFA) layer sits in the middle. When I_O falls on a color image sensor, photons first go through the MLA to improve light collection. Next, light passes through the CFA layer, resulting in a mosaic pattern of the three stimulus RGB colors. Finally, the bottom layer collects the color-filtered light and outputs a single channel RAW image, I_W .

The detailed process is illustrated in Figure 2. In particular, at the bottom layer, a potential well counts photons arriving at its location (x, y) and converts the accumulated photons into electrons, and the conversion process is specified by the detector quantum efficiency. During the process, electrons could be generated by other resources, called electron noise. Two common electron noise types are the dark noise n_d , which is independent of light; and dark current n_I , which depends on the sensor temperature. These follow normal and Poisson distributions, respectively [75]. Next, the converted electrons are clipped based on the maximum well capacity, e_s , and scaled by a sensor gain factor κ . Finally, the modulated electrons are converted to digital values by an analog-to-digital converter (ADC), which involves quantization of the input and introduces a small amount of noise, n_q .

Mathematically, a pixel of a RAW image, I_W , at position (x, y) can be defined as:

$$I_W(x, y) = b + n_q + \kappa \min(e_s, n_d + n_I + \sum_{\lambda} e(x, y, \lambda)), \quad (1)$$

where b is the black level, level of brightness with no light; $e(x, y, \lambda)$ is the number of electrons arrived at a well at position (x, y) for wavelength λ .

This image formation model reveals that besides the natural scene being captured, RAW images heavily depend on the *specific stochastic natures of the optics, color filtering, sensing, and readout components*. The proposed method exploits these statistics.

4 Raw Image Domain Defense

In this section, we describe the proposed defense method, which leverages the distribution of RAW measurements as a prior to project adversarially perturbed RGB images to benign ones. Given an adversarial input, existing defense approaches learn an RGB-to-RGB projection from the adversarially perturbed distribution of RGB images, Y' , to the closest point in corresponding RGB natural distribution, Y . We use the operator $T : Y' \rightarrow Y$ for this projection operation. As this RGB distribution Y empirically sampled from the ISP outputs of diverse existing cameras, it also ingests diverse reconstruction artifacts, making it impossible to exploit photon-flux specific cues, e.g., photon shot noise, optical aberrations, or camera-specific readout characteristics – as image processing pipelines are designed to remove such RAW cues.

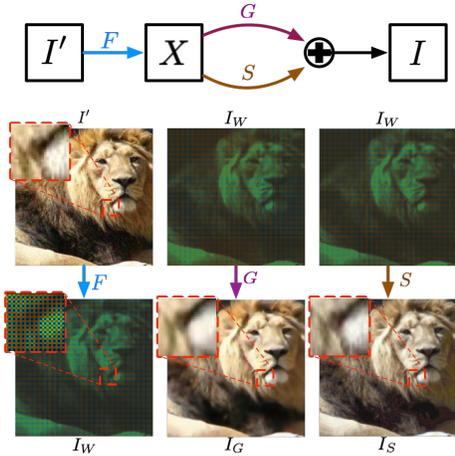


Fig. 3: Overview of the proposed defense approach, see text. We note that the resolution of the RAW image (RGGGB) is twice larger than that of the RGB image. We linearly scaled the RAW image in this figure for better visualization.

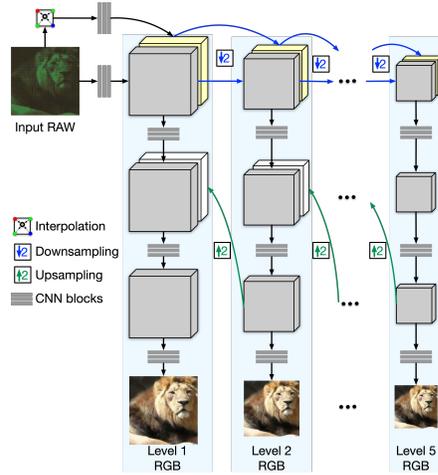


Fig. 4: The architecture of the G operator, which is adopted and modified from PyNet [33]. The finer operator level exploits upsampled coarser-level features to reconstruct the RGB output. The model is trained sequentially in a coarse-to-fine manner.

Departing from existing methods, as illustrated in Figure 1, 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, multi-spectral 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. Operators F and G are trained separately without end-to-end fine-tuning. Notably, the proposed defense scheme is *entirely model-agnostic* as it does not require any knowledge of potential adversarial attacks.

For defending against an attack, as shown in Figure 3, the proposed approach first uses the F operator to map an input adversarial image, I' , to its intermediate RAW measurements, I_W . Then, I_W is processed separately by the G and S operators to convert it to two images in the natural RGB distribution, denoted

as I_G and I_S , 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. Mathematically, the defense process is defined as:

$$I = \omega G(F(I')) + (1 - \omega)S(F(I')), \quad (2)$$

where ω is a hyper-parameter for weighting the contributions from the two operators G and S . In the following sections, we introduce each operator in detail.

4.1 F Operator: Image to RAW Mapping

We use a small learned encoder-decoder network as the F operator to map an RGB image to its intermediate RAW measurements. The details of network architecture is shown in supplementary.

We train this module in a supervised manner with two \mathcal{L}_2 losses. Both of the \mathcal{L}_2 losses are calculated based on ground truth (GT) RAW and estimated RAW images. The only difference between the two losses is the input RGB image used for evaluating a RAW image. One is with the original input RGB image, while the other is generated by adding Gaussian noise to the original input RGB image. In doing so, F is trained with the ability to convert both benign and slightly perturbed RGB images to their RAW distribution. We note that the added Gaussian distribution is *different from the correlated noise generated by various adversarial attacks*. Mathematically, given a benign RGB image, I , and its corresponding GT RAW measurements, GT_W , the loss function is defined as:

$$\mathcal{L}_F = \|F(I), GT_W\|_2 + \|F(I + \alpha\varepsilon), GT_W\|_2, \quad (3)$$

$$\varepsilon \sim \mathcal{N}(\mu, \sigma), \quad (4)$$

where ε is a Gaussian noise with mean, μ , and standard deviation, σ ; α is a random number in the range between 0 and 1, weighting the amount of noise added. We empirically set the μ and σ to 0 and 1, respectively.

4.2 G Operator: Learned ISP

The G operator, learned network, converts the I_w generated by the F operator to an RGB image. The challenge of converting RAW images to RGB images is that the process requires both global and local modifications. The global modifications aim to change the high-level properties of the image, such as brightness and white balance. In contrast, the local modifications refer to low-level processing like texture enhancement, sharpening, and noise removal. During the image reconstruction process, an effective local process is expected to be guided by global contextual information, which requires the information exchange between global and local operations. This motivates us to leverage a pyramidal convolutional neural network to fuse global and local features for optimal reconstruction results. We adopt and modify architecture similar to the PyNet [33]. As shown in Figure 4, the network has five levels, the 1st level is the finest, and the 5th level is the coarsest. The finer-level uses upsampled features from the coarser-level by concatenating them. We modified PyNet by adding an interpolation

layer before the input of each level, interpolating the downsampled RAW Bayer pattern. This practice facilitates learning as the network only needs to learn the residuals between interpolated RGB and ground truth RGB, leading to better model performance.

The loss function for this model consists of three components: perceptual, structural similarity, and \mathcal{L}_2 loss. The perceptual and \mathcal{L}_2 loss functions are adopted to ensure the fidelity of the reconstructed image, and the structural similarity loss function [81] is used to enhance the dynamic range. Given an input RAW image, I_W , and the corresponding GT RGB image GT_I , the loss function can be mathematically defined as:

$$\begin{aligned} \mathcal{L}_G^i = & \beta^i \mathcal{L}_{Perc}(G(I_W), GT_I) + \gamma^i \mathcal{L}_{SSIM}((G(I_W), GT_I)) \\ & + \mathcal{L}_2(G(I_W), GT_I) \quad \text{for } i \in [1, 5], \end{aligned} \quad (5)$$

where i represents the training level. As the model is trained in a coarse-to-fine manner, different losses are used for each level i . \mathcal{L}_{Perc} , \mathcal{L}_{SSIM} , and \mathcal{L}_2 represents the perceptual loss calculated with VGG architecture, structural similarity loss, and \mathcal{L}_2 loss, respectively; β^i and γ^i are the two weighting hyper-parameters, which are set empirically. The model is trained sequentially in a coarse-to-fine manner, *i.e.*, from $i = 5$ to $i = 1$.

4.3 S Operator: Conventional ISP

The S operator has the same functionality as the G operator, converting a RAW image to an RGB image. Unlike the G operator, the S operator offers the functionalities of a conventional hardware ISP pipeline using a sequence of cascaded sub-modules, and it is non-differentiable.

While we may exploit the ISP pipeline of any digital camera we can extract raw and post-ISP data from, we use a software-based ISP pipeline consisting of the following components: Bayer demosaicing, color balancing, white balancing, contrast improvement, and colorspace conversion sub-modules. Based on the Zurich-Raw-to-RGB dataset [34], we manually tune the hyperparameters of all sub-modules to find the optimal ones that offer the converted RGB image with similar image quality to the original RGB ones. We refer the reader to the Supplementary Material for a detailed description.

4.4 Operator Training

We use the Zurich-Raw-to-RGB dataset [34] to train the F and G operators. The Zurich-Raw-to-RGB dataset consists of 20,000 RAW-RGB image pairs, captured using a Huawei P20 smartphone with a 12.3 MP Sony Exmor IMX380 sensor and a Canon 5D Mark IV DSLR. Both of the F and G operators are trained in PyTorch with Adam optimizer on NVIDIA A100 GPUs. We set the learning rate to $1e-4$ and $5e-5$ for training F and G operators, respectively. The hyperparameters used in our approach have the following settings: $\omega = 0.7$ in Eq. 2; $\mu = 0$ and $\sigma = 1$ for the Gaussian noise used in Eq. 4; In Eq. 5, β^i is set to 1 for $i \in [1, 3]$ and 0 for $i \in [4, 5]$; γ^i is set to 1 for $i = 1$ and 0 for $i \in [2, 5]$.

	FSGM		PGD		BIM		DeepFool		C&W		NewtonFool BPDA	
	2/255 \uparrow	4/255 \uparrow	2/255 \uparrow	4/255 \uparrow	2/255 \uparrow	4/255 \uparrow	$L_\infty \uparrow$	$L_2 \uparrow$	$L_\infty \uparrow$	$L_2 \uparrow$	$L_\infty \uparrow$	$L_\infty \uparrow$
ResNet-101												
JPEG-Defense [20]	33.14	20.71	45.19	21.74	36.78	8.5	53.16	45.69	59.06	52.01	24.65	0.08
TVM [27]	43.75	40.02	45.46	44.35	44.86	41.93	47.69	39.89	45.51	40.44	22.6	6.39
Randomized Resizing & Padding [84]	45.21	34.97	45.38	27.75	40.04	18.04	73.06	62.47	66.53	59.87	27.93	2.66
HGD [47]	54.75	43.85	55.26	50.05	56.74	48.61	64.34	58.13	59.98	52.88	27.70	0.03
Pixel-Deflection [64]	54.56	35.14	60.68	34.86	58.71	41.91	75.97	64.13	66.29	60.91	28.81	1.87
ComDefend [38]	48.21	36.51	53.28	48.38	51.39	42.01	63.68	55.62	58.53	50.38	26.46	0.03
Proposed Method	66.02	58.85	68.34	66.17	66.91	63.01	72.04	63.52	71.40	67.33	40.96	38.85
InceptionV3												
JPEG-Defense [20]	31.97	20.25	43.34	21.15	34.68	8.55	51.20	43.49	55.00	50.39	24.06	0.12
TVM [27]	42.47	37.23	42.75	41.61	42.80	39.71	45.21	37.39	43.27	37.51	23.05	4.58
Randomized Resizing & Padding [84]	41.86	34.49	43.41	25.60	39.42	16.62	70.24	58.65	63.24	55.62	27.55	2.09
HGD [47]	52.83	40.99	50.35	47.62	56.02	47.78	60.33	56.61	59.55	52.0	26.84	0.03
Pixel-Deflection [64]	51.42	34.27	56.13	32.49	56.18	39.13	71.16	61.58	61.94	57.58	28.01	1.56
ComDefend [38]	47.00	35.34	49.99	46.15	48.74	39.58	60.01	52.47	55.85	47.70	25.44	0.03
Proposed Method	63.03	56.34	65.69	63.03	64.77	59.49	69.25	60.04	66.97	64.69	38.01	36.43

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.

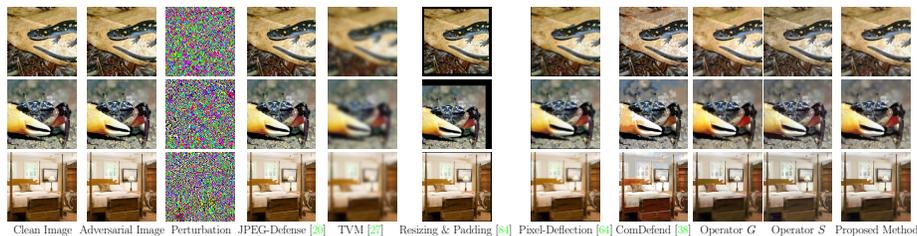


Fig. 5: Qualitative outputs of the proposed method along with both G and S operators, and state-of-the-art defense methods on the ImageNet dataset, see text.

5 Experiments & Analysis

The proposed method acts as an off-the-shelf input preprocessing module, and it requires no additional training to be transferred to different tasks. To validate the effectiveness and generalization capabilities of the proposed defense approach, we evaluate the method on three different vision tasks, *i.e.*, classification, semantic segmentation, and 2D object detection, with corresponding adversarial attacks.

5.1 Experimental Setup

Adversarial Attack Methods: We evaluate our method by defending against the following attacks: FGSM [26], BIM [43], PGD [53], C&W [8], NewtonFool [37], and DeepFool [54]. For classification, we use the widely used Foolbox benchmarking suite [65] to implement these attack methods. Since Foolbox does not directly support semantic segmentation and object detection, we use the lightweight TorchAttacks library [41] for generating adversarial examples with FGSM, PGD, and BIM attacks. We also evaluate against the DAG [85] attack, a dedicated attack approach for semantic segmentation and object detection tasks. Moreover, we further evaluate against BPDA [2], an attack method specifically designed for circumventing input transformation defenses that rely on obfuscated gradients. Applying our method for defending against BPDA, however, requires a slight modification at inference time, see Supplementary Document for details. We note that all applied attacks are untargeted. Definitions of all attack methods are provided in the Supplementary Material.

	FSGM		PGD		BIM		DAG	
	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$
JPEG-Defense [20]	37.41	32.27	24.53	6.21	25.74	10.18	14.12	5.66
TVM [27] [27]	42.64	41.53	45.55	42.24	44.51	38.44	31.88	25.56
HGD [47]	43.39	40.82	44.54	40.88	40.03	39.95	28.61	22.36
Pixel-Deflection [64]	44.13	41.88	46.38	42.32	44.78	37.22	30.73	24.61
ComDefend [38]	45.57	39.23	44.85	41.14	42.71	36.12	28.94	23.36
Proposed Method	52.35	48.04	53.41	49.59	54.86	50.51	40.35	37.88

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.

	FSGM		PGD		BIM		DAG	
	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$	$L_\infty = 2/255 \uparrow$	$L_\infty = 4/255 \uparrow$
JPEG-Defense [20]	39.02	35.88	37.96	33.51	38.85	34.69	30.72	25.07
TVM [27]	48.11	39.66	47.1	44.38	48.94	41.76	39.20	33.18
HGD [47]	50.68	40.06	51.24	45.92	46.80	39.74	41.15	37.23
Pixel-Deflection [64]	53.77	44.82	54.45	47.22	55.32	48.32	46.52	39.87
ComDefend [38]	50.18	42.93	50.46	43.08	52.32	44.2	44.68	37.22
Proposed Method	61.68	59.37	64.71	60.23	66.52	61.82	57.83	54.12

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.

Baseline Defense Approaches: We compare to the following input transformation defense methods: JPEG compression [20], randomized resizing & padding [84], image quilting [27], TVM [27], HGD [47], pixel deflection [64], and Comdefend [38]. We evaluate all baseline methods on the three vision tasks, except that the randomized resizing & padding method is omitted in semantic segmentation and object detection tasks as it destroys the semantic structure. We directly adopt the open-source PyTorch implementation for all baseline methods. We use the same training dataset as the one used to train our method for those methods that required training. It is worth noting that all baseline methods do not require adversarial examples for training.

Evaluation Dataset and Metrics: For classification, we use the ImageNet validation set and evaluate the Top-1 classification accuracy of all competing defense approaches. For semantic segmentation and object detection, we evaluate on the MS COCO [48] and Pascal VOC [21] datasets. The effectiveness of all methods for segmentation and detection is measured by mean Intersection over Union (mIoU) and mean Average Precision (mAP), respectively.

5.2 Assessment

Classification: We apply a given attack method with ResNet101 and InceptionV3 to generate adversarial samples. For FSGM, BIM, and PGD, we set two different maximum perturbation levels in L_∞ distance, namely $2/255$ and $4/255$. The maximum number of iterations is set to 100 for both BIM and PGD. For C&W, NewtonFool, and DeepFool attacks, we generate both L_∞ distance based attacks and L_2 distance-based attacks; we choose 100 update steps for C&W and NewtonFool, and 50 for DeepFool; DeepFool requires the number of candidate classes, which is set to 10 in our experiments.

The Top-1 classification accuracies of all methods are reported in Table 1. Our approach outperforms the baseline methods with a large margin under all experimental settings except those with DeepFool attacks. Notably, under DeepFool attacks, the differences between the best performer pixel-deflection and ours

are marginal. Moreover, with PGD and BIM attacks, our defense method offers the lowest relative performance degradation when the more vigorous attack is performed (*i.e.*, maximum perturbation increases from 2/255 to 4/255). Figure 5 qualitatively underlines the motivation of combining G and S operators in a weighted sum manner. The G operator learns to mitigate the adversarial pattern, *i.e.*, it recovers a latent image in the presence of severe measurement uncertainty, while the S operator can faithfully reconstruct high-frequency details. Note that our method is able to generalize well to images from the ImageNet dataset, which typically depict single objects, although it is trained on the Zurich-Raw-to-RGB dataset, consisting of street scenes.

Semantic Segmentation: In this task, we conduct experiments with two different types of attacks: the commonly used adversarial attacks, and the attack specially designed for attacking semantic segmentation models. For the former, FGSM, BIM, and PGD are used; We use DAG [85], a dedicated semantic segmentation attack, for the latter. All attacks are based on a COCO-pretrained DeepLabV3 model [12]. Two different maximum perturbation levels in L_∞ are used (*i.e.*, 2/255 and 4/255). The corresponding experimental results are reported in Table 2. The proposed approach significantly outperforms baseline methods under all experimental settings. Note that no additional training is required to apply the proposed Raw-Defense approach to defend other vision tasks, validating the generalization capabilities of the method.

2D Object Detection: The experimental settings are the same as the ones used for semantic segmentation experiments, except that we use a pretrained Faster R-CNN [66]. We report the mAPs on the Pascal VOC dataset under different experimental settings in Table 3. The proposed defense method offers the best defense performance in all experimental settings, indicating that our approach generalizes well to unseen tasks.

5.3 RAW Distribution Analysis

In this section, we provide additional analysis on the function of the RAW distribution as an intermediate mapping space. Fundamentally, we share the motivation from existing work that successfully exploits RAW data for imaging and vision tasks, including [17, 77]. RGB images are generated by processing RAW sensor measurements (see Sec. 3) with an image processing pipeline. This process removes statistical information embedded in the sensor measurements by aberrations in the optics, readout noise, color filtering, exposure, and scene illumination. While existing work directly uses RAW inputs to preserve this information, we exploit it in the form of an *empirical intermediate image distribution*. Specifically, we devise a mapping via RAW space, thereby using RAW data to train network mapping modules, which we validate further below. As a result, we allow the method to remove adversarial patterns not only by relying on RGB image priors but also RAW image priors. We *validate the role of RAW data* in our method in Table 4, resulting in a large Top-1 accuracy drop (*i.e.*, more than 12%), when swapping the real RAW distribution to a synthesized one. This is further corroborated in Table 5, where the defense breaks down from 71% to 53%, when gradually moving from RAW to RGB as intermediate image space.

These experiment validate that, the “rawer” the intermediate image space is, the better the defense performs.

Effect of Intermediate Mapping Space: We use the RAW image distribution as the intermediate mapping space in our method. To validate the effectiveness of this choice, we map to other intermediate stages in the processing pipeline, such as demosaicing stage, color balance stage, and the white balance stage. Specifically, we assess how using different stage values as intermediate mapping space affects the defense performance (*i.e.*, we ablate on the intermediate mapping space used). As reported in Table 5, we observe that the defense performance gradually decreases as we map via a less RAW intermediate space. In other words, *the “rawer” the intermediate image space is, the better performance can be achieved.* This validates the importance and benefit of exploiting RAW distribution in the defense.

Real RAW Versus Synthetic RAW: We further ablate on the dataset used to train our model. Specifically, we trained F and G operator on the Zurich-Raw-to-RGB dataset, HDR-RAW-RGB [29] and MIT-RAW-RGB [23] respectively and assess how the defense performance changes. Similar to Zurich-Raw-to-RGB dataset, the RAW images in HDR-RAW-RGB are captured by a real camera; however, the ones offered by MIT-RAW-RGB are purely synthesized by reformatting downsampled RGB images into Bayer patterns with handcrafted Gaussian noise. We note that, as such, the MIT-RAW-RGB dataset does not include the RAW distribution cues. The experimental results are reported in Table 4. As observed, both RAW distribution Zurich-Raw-to-RGB and HDR-RAW-RGB allow us to learn effective adversarial defenses, while a sharp performance degradation occurs when shifting from real RAW distribution to the synthesized one due to the lack of natural RAW distribution cues. This again validates the effectiveness of statistical information in the real RAW distribution when defending against adversarial attacks.

	FSGM	PGD	C&W	NewtonFool	DeepFool
Zurich-Raw-to-RGB [34]	58.85	66.17	71.40	40.96	72.04
HDR-RAW-RGB [29]	55.57	64.12	71.65	42.36	70.77
MIT-RAW-RGB [23]	40.52	47.29	55.13	28.52	58.49

Table 4: **Quantitative Ablation Study on RAW Training Datasets.** We train F and G two operators with three different RAW-RGB datasets and report the Top-1 defense accuracy on the ImageNet dataset. The RAW images in the Zurich-Raw-to-RGB and HDR-RAW-RGB are captured by real cameras, while the ones in MIT-RAW-RGB are synthesized. We see a sharp performance drop when swapping the real RAW training data to synthetic data due to the lack of natural RAW distribution cues.

	FSGM	PGD	C&W	NewtonFool	DeepFool
Raw Capture	57.33	65.02	70.86	40.65	70.23
Demosaic Stage	52.93	59.83	63.29	36.76	64.81
Color Balance Stage	48.38	55.41	57.92	33.92	59.37
White Balance Stage	47.35	54.02	56.23	33.01	57.08
contrast Improvement Stage	45.2	52.18	54.18	31.84	55.64
Agamma adjustment Stage	44.4	50.91	53.08	30.57	54.19

Table 5: **Effect of Different Intermediate Mapping Spaces.** We report the Top-1 adversarial defense accuracy on ImageNet dataset when mapping to different intermediate mapping spaces that are the steps of the image processing pipeline. The performance drops as the intermediate image space moves from RAW to the RGB output space. This validates the importance and benefit of exploiting RAW distribution in the defense.

5.4 Robustness to Hyper-parameter and Operator Deviations

Hyper-parameter ω : We introduced a hyper-parameter ω for weighting the contributions of the two operators G and S . Next, we evaluate how varying values of ω affect the overall defense accuracy. As reported in Tab. 6, we find that, while

each attack has a different optimal value of ω , the range 0.6-0.8 provides a good trade-off, and we use 0.7 in our experiments. Limited by space, we report here a subset of attacks.

Hyper-parameter $\omega =$	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
Against FSGM Attack	64.25	64.41	64.83	65.27	65.58	65.87	65.93	66.02	65.75	65.53	65.39
Against C&W Attack	69.16	69.44	69.93	70.26	70.81	70.96	71.35	71.40	72.70	71.28	71.07
Against DeepFool Attack	69.55	69.84	71.19	71.51	71.88	72.35	72.63	72.04	71.75	71.69	71.04

Table 6: **Effect of hyper-parameter ω .** We evaluate the impact of the method hyper-parameter ω on the effectiveness of the proposed defense method.

	F -300	F -320	F -340	F -360	F -380	F -400	Gaussian Noise σ	0 (no noise)	0.01	0.05	0.1	0.3	0.5
G -300	66.02	66.08	65.93	66.05	66.03	66.05	Against FSGM Attack	66.02	66.01	65.98	65.90	65.73	65.64
G -330	66.11	66.04	65.97	66.01	65.99	66.04	Against PGD Attack	68.34	68.30	68.22	68.10	68.03	67.83
G -360	65.98	66.04	66.02	65.89	66.08	69.01							
G -390	65.98	65.95	66.00	66.04	65.99	65.94							

Table 7: **Robustness to Deviations of F and G .** We evaluate the defense accuracy when mixing operator from different training epochs.

Table 8: **Robustness to Deviations of F .** We perturb the output of operator F with Gaussian noise of different standard deviations and report the defense accuracy.

Deviations of Operators F and G : The operator F and G are trained separately and used jointly at the inference time. We evaluate how deviations in each operator affect the overall performance in two experiments. First, we mix the operators F and G from different training checkpoints and evaluate the effect on the defense accuracy. Tab. 7 reports that the checkpoint combinations do not result in a failure but only slight deviations of the defense performance. Second, we add varying levels of Gaussian noise $G(0, \sigma)$ to the output of operator F and evaluate how such deviation affects the following steps and the overall defense accuracy. Tab. 8 reports that such perturbations are not amplified in the following steps, and the defense accuracy only fluctuates slightly. The experiments show that the ISP operators G and S themselves are robust to slight deviation in each component.

6 Conclusion

We exploit RAW image data as an empirical latent space in the formulation of the proposed adversarial defense method. Departing from existing defense methods that aim to directly map an adversarially perturbed image to the closest benign image, we exploit large-scale natural image datasets as an empirical prior for sensor-captured images – before they end up in existing datasets after their transformation through conventional image processing pipelines. This empirical prior allows us to rely on low-level image processing pipelines to design the mappings between the benign and perturbed image distributions. We validate the effectiveness of the method, which is entirely model-agnostic, requires no adversarial examples to train, and acts as an off-the-shelf preprocessing module that can be transferred to diverse tasks. We also provided insight into the working principles of the approach and assess that the method significantly outperforms the comparable baselines. In the future, we plan to explore RAW natural image statistics as an unsupervised prior for image reconstruction and generative neural rendering tasks.

References

1. Akhtar, N., Mian, A., Kardan, N., Shah, M.: Advances in adversarial attacks and defenses in computer vision: A survey (2021) [4](#)
2. Athalye, A., Carlini, N., Wagner, D.: Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In: International conference on machine learning. pp. 274–283. PMLR (2018) [10](#)
3. Athalye, A., Engstrom, L., Ilyas, A., Kwok, K.: Synthesizing robust adversarial examples. In: International conference on machine learning. pp. 284–293. PMLR (2018) [4](#)
4. Bahat, Y., Irani, M., Shakhnarovich, G.: Natural and adversarial error detection using invariance to image transformations. arXiv preprint arXiv:1902.00236 (2019) [4](#)
5. Borkar, T., Heide, F., Karam, L.: Defending against universal attacks through selective feature regeneration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 709–719 (2020) [1](#)
6. Brendel, W., Rauber, J., Bethge, M.: Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In: International Conference on Learning Representations (2018) [4](#)
7. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE Symposium on Security and Privacy (2017) [4](#)
8. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks (2017) [10](#)
9. Chen, C., Chen, Q., Do, M.N., Koltun, V.: Seeing motion in the dark. In: 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019. pp. 3184–3193. IEEE (2019). <https://doi.org/10.1109/ICCV.2019.00328>, <https://doi.org/10.1109/ICCV.2019.00328> [3](#)
10. Chen, C., Chen, Q., Xu, J., Koltun, V.: Learning to see in the dark. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18–22, 2018. pp. 3291–3300. Computer Vision Foundation / IEEE Computer Society (2018). <https://doi.org/10.1109/CVPR.2018.00347>, http://openaccess.thecvf.com/content_cvpr_2018/html/Chen_Learning_to_See_CVPR_2018_paper.html [3](#)
11. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence **40**(4), 834–848 (2017) [1](#)
12. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation (2017) [12](#)
13. Chen, P.Y., Zhang, H., Sharma, Y., Yi, J., Hsieh, C.J.: Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In: Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security. pp. 15–26 (2017) [4](#)
14. Cheng, M., Le, T., Chen, P.Y., Yi, J., Zhang, H., Hsieh, C.J.: Query-efficient hard-label black-box attack: An optimization-based approach. arXiv preprint arXiv:1807.04457 (2018) [4](#)
15. Dai, L., Liu, X., Li, C., Chen, J.: Awnet: Attentive wavelet network for image isp. In: ECCV Workshops (2020) [3](#)

16. Das, N., Shanbhogue, M., Chen, S.T., Hohman, F., Chen, L., Kounavis, M.E., Chau, D.H.: Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression. arXiv preprint arXiv:1705.02900 (2017) [2](#), [5](#)
17. Diamond, S., Sitzmann, V., Julca-Aguilar, F., Boyd, S., Wetzstein, G., Heide, F.: Dirty pixels: Towards end-to-end image processing and perception. ACM Transactions on Graphics (SIGGRAPH) (2021) [12](#)
18. Duan, R., Ma, X., Wang, Y., Bailey, J., Qin, A.K., Yang, Y.: Adversarial camouflage: Hiding physical-world attacks with natural styles. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1000–1008 (2020) [4](#)
19. Dziugaite, G.K., Ghahramani, Z., Roy, D.M.: A study of the effect of jpg compression on adversarial images (2016) [2](#)
20. Dziugaite, G.K., Ghahramani, Z., Roy, D.M.: A study of the effect of JPG compression on adversarial images. CoRR [abs/1608.00853](#) (2016) [4](#), [5](#), [10](#), [11](#)
21. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International journal of computer vision **88**(2), 303–338 (2010) [11](#)
22. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., Song, D.: Robust physical-world attacks on deep learning visual classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1625–1634 (2018) [4](#)
23. Gharbi, M., Chaurasia, G., Paris, S., Durand, F.: Deep joint demosaicking and denoising. ACM Transactions on Graphics (TOG) **35**(6), 191 (2016) [13](#)
24. Gong, C., Ren, T., Ye, M., Liu, Q.: Maxup: Lightweight adversarial training with data augmentation improves neural network training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2474–2483 (June 2021) [2](#), [5](#)
25. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. CoRR [abs/1412.6572](#) (2015) [4](#)
26. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples (2015) [10](#)
27. Guo, C., Rana, M., Cisse, M., Van Der Maaten, L.: Countering adversarial images using input transformations. ICLR (2018) [2](#), [5](#), [10](#), [11](#)
28. Guo, M., Yang, Y., Xu, R., Liu, Z., Lin, D.: When nas meets robustness: In search of robust architectures against adversarial attacks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 631–640 (2020) [4](#)
29. Hasinoff, S.W., Sharlet, D., Geiss, R., Adams, A., Barron, J.T., Kainz, F., Chen, J., Levoy, M.: Burst photography for high dynamic range and low-light imaging on mobile cameras. ACM Transactions on Graphics (ToG) **35**(6), 1–12 (2016) [13](#)
30. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017) [1](#)
31. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) [1](#)
32. Hu, W., Tan, Y.: Generating adversarial malware examples for black-box attacks based on gan. ArXiv [abs/1702.05983](#) (2017) [4](#)
33. Ignatov, A., Gool, L.V., Timofte, R.: Replacing mobile camera isp with a single deep learning model (2020) [7](#), [8](#)
34. Ignatov, A.D., Gool, L.V., Timofte, R.: Replacing mobile camera isp with a single deep learning model. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) pp. 2275–2285 (2020) [3](#), [9](#), [13](#)

35. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017) [1](#)
36. Jan, S.T., Messou, J., Lin, Y.C., Huang, J.B., Wang, G.: Connecting the digital and physical world: Improving the robustness of adversarial attacks. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 962–969 (2019) [4](#)
37. Jang, U., Wu, X., Jha, S.: Objective metrics and gradient descent algorithms for adversarial examples in machine learning. In: Proceedings of the 33rd Annual Computer Security Applications Conference. pp. 262–277 (2017) [10](#)
38. Jia, X., Wei, X., Cao, X., Foroosh, H.: Comdefend: An efficient image compression model to defend adversarial examples. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6084–6092 (2019) [2](#), [5](#), [10](#), [11](#)
39. Karaimer, H.C., Brown, M.S.: A software platform for manipulating the camera imaging pipeline. In: ECCV (2016) [3](#)
40. Karaimer, H.C., Brown, M.S.: A software platform for manipulating the camera imaging pipeline. In: European Conference on Computer Vision. pp. 429–444. Springer (2016) [5](#)
41. Kim, H.: Torchattacks: A pytorch repository for adversarial attacks (2021) [10](#)
42. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world. arXiv preprint arXiv:1607.02533 (2016) [4](#)
43. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world (2017) [10](#)
44. Li, Y., Li, L., Wang, L., Zhang, T., Gong, B.: Nattack: Learning the distributions of adversarial examples for an improved black-box attack on deep neural networks. arXiv preprint arXiv:1905.00441 (2019) [4](#)
45. Liang, Z., Cai, J., Cao, Z., Zhang, L.: Cameranet: A two-stage framework for effective camera isp learning. IEEE Transactions on Image Processing **30**, 2248–2262 (2021). <https://doi.org/10.1109/TIP.2021.3051486> [3](#)
46. Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., Zhu, J.: Defense against adversarial attacks using high-level representation guided denoiser. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1778–1787 (2018) [2](#), [5](#)
47. Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., Zhu, J.: Defense against adversarial attacks using high-level representation guided denoiser. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1778–1787 (2018) [10](#), [11](#)
48. Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft coco: Common objects in context (2015) [11](#)
49. Liu, Z., Liu, Q., Liu, T., Wang, Y., Wen, W.: Feature Distillation: DNN-oriented JPEG compression against adversarial examples. International Joint Conference on Artificial Intelligence (2018) [5](#)
50. Liu, Z., Liu, Q., Liu, T., Xu, N., Lin, X., Wang, Y., Wen, W.: Feature distillation: Dnn-oriented jpeg compression against adversarial examples. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 860–868. IEEE (2019) [2](#), [5](#)
51. Lu, J., Issaranon, T., Forsyth, D.: Safetynet: Detecting and rejecting adversarial examples robustly. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 446–454 (2017) [1](#)

52. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017) **1, 4**
53. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks (2019) **10**
54. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks (2016) **10**
55. Mosleh, A., Sharma, A., Onzon, E., Mannan, F., Robidoux, N., Heide, F.: Hardware-in-the-loop end-to-end optimization of camera image processing pipelines. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) **3, 4**
56. Nakkiran, P.: Adversarial robustness may be at odds with simplicity. arXiv preprint arXiv:1901.00532 (2019) **1**
57. Narodytska, N., Kasiviswanathan, S.: Simple black-box adversarial attacks on deep neural networks. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 1310–1318 (2017). <https://doi.org/10.1109/CVPRW.2017.172> **4**
58. Pang, T., Xu, K., Dong, Y., Du, C., Chen, N., Zhu, J.: Rethinking softmax cross-entropy loss for adversarial robustness. ICLR (2020) **4**
59. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against machine learning. In: Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security. p. 506–519. ASIA CCS '17, Association for Computing Machinery, New York, NY, USA (2017). <https://doi.org/10.1145/3052973.3053009>, <https://doi.org/10.1145/3052973.3053009> **4**
60. Papernot, N., McDaniel, P., Wu, X., Jha, S., Swami, A.: Distillation as a defense to adversarial perturbations against deep neural networks. In: 2016 IEEE Symposium on Security and Privacy (SP). pp. 582–597. IEEE (2016) **1, 4**
61. Papernot, N., McDaniel, P.D., Goodfellow, I.J.: Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. CoRR **abs/1605.07277** (2016), <http://arxiv.org/abs/1605.07277> **4**
62. Phan, B., Mannan, F., Heide, F.: Adversarial imaging pipelines. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16051–16061 (2021) **4**
63. Poursaeed, O., Katsman, I., Gao, B., Belongie, S.: Generative adversarial perturbations. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4422–4431 (2018) **1**
64. Prakash, A., Moran, N., Garber, S., DiLillo, A., Storer, J.: Deflecting adversarial attacks with pixel deflection. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (2018) **5, 10, 11**
65. Rauber, J., Brendel, W., Bethge, M.: Foolbox: A python toolbox to benchmark the robustness of machine learning models (2018) **10**
66. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks (2016) **12**
67. Samangouei, P., Kabkab, M., Chellappa, R.: Defense-gan: Protecting classifiers against adversarial attacks using generative models. ICLR (2018) **2, 5**
68. Schwartz, E., Giryes, R., Bronstein, A.M.: Deepisp: Toward learning an end-to-end image processing pipeline **28(2)**, 912–923 (Feb 2019). <https://doi.org/10.1109/TIP.2018.2872858>, <https://doi.org/10.1109/TIP.2018.2872858> **3**

69. Sen, S., Ravindran, B., Raghunathan, A.: Empir: Ensembles of mixed precision deep networks for increased robustness against adversarial attacks. *ICLR (2020)* [4](#)
70. Shafahi, A., Najibi, M., Ghiasi, A., Xu, Z., Dickerson, J., Studer, C., Davis, L.S., Taylor, G., Goldstein, T.: Adversarial training for free! In: *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. pp. 3358–3369 (2019) [5](#)
71. Shi, Y., Wang, S., Han, Y.: Curls & whey: Boosting black-box adversarial attacks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 6512–6520 (2019) [4](#)
72. Stutz, D., Hein, M., Schiele, B.: Confidence-calibrated adversarial training: Generalizing to unseen attacks. In: *International Conference on Machine Learning*. pp. 9155–9166. *PMLR (2020)* [5](#)
73. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199 (2013)* [1](#), [4](#)
74. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks (2014) [1](#)
75. Tseng, E., Mosleh, A., Mannan, F., St-Arnaud, K., Sharma, A., Peng, Y., Braun, A., Nowrouzezahrai, D., Lalonde, J.F., Heide, F.: Differentiable compound optics and processing pipeline optimization for end-to-end camera design. *ACM Transactions on Graphics (TOG)* **40**(4) (2021) [6](#)
76. Tseng, E., Yu, F., Yang, Y., Mannan, F., Arnaud, K.S., Nowrouzezahrai, D., Lalonde, J.F., Heide, F.: Hyperparameter optimization in black-box image processing using differentiable proxies **38**(4) (Jul 2019). <https://doi.org/10.1145/3306346.3322996>, <https://doi.org/10.1145/3306346.3322996> [3](#), [4](#)
77. Tseng, E., Yu, F., Yang, Y., Mannan, F., Arnaud, K.S., Nowrouzezahrai, D., Lalonde, J.F., Heide, F.: Hyperparameter optimization in black-box image processing using differentiable proxies. *ACM Trans. Graph.* **38**(4), 27–1 (2019) [12](#)
78. Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., Madry, A.: Robustness may be at odds with accuracy. In: *International Conference on Learning Representations*. No. 2019 (2019) [1](#)
79. Tu, C.C., Ting, P., Chen, P.Y., Liu, S., Zhang, H., Yi, J., Hsieh, C.J., Cheng, S.M.: Autozoo: Autoencoder-based zeroth order optimization method for attacking black-box neural networks. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 33, pp. 742–749 (2019) [4](#)
80. Wang, J., Zhang, H.: Bilateral adversarial training: Towards fast training of more robust models against adversarial attacks. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 6629–6638 (2019) [2](#), [5](#)
81. Wang, Z., Simoncelli, E., Bovik, A.: Multiscale structural similarity for image quality assessment. *The Thirty-Seventh Asilomar Conference on Signals, Systems and Computers, 2003* . <https://doi.org/10.1109/acssc.2003.1292216> [9](#)
82. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training. *ICLR (2020)* [2](#), [5](#)
83. Wu, Y.H., Yuan, C.H., Wu, S.H.: Adversarial robustness via runtime masking and cleansing. In: *International Conference on Machine Learning*. pp. 10399–10409. *PMLR (2020)* [4](#)
84. Xie, C., Wang, J., Zhang, Z., Ren, Z., Yuille, A.: Mitigating adversarial effects through randomization. *arXiv preprint arXiv:1711.01991 (2017)* [5](#), [10](#), [11](#)
85. Xie, C., Wang, J., Zhang, Z., Zhou, Y., Xie, L., Yuille, A.: Adversarial examples for semantic segmentation and object detection (2017) [10](#), [12](#)

86. Xie, C., Wu, Y., Maaten, L.v.d., Yuille, A.L., He, K.: Feature denoising for improving adversarial robustness. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 501–509 (2019) [1](#)
87. Xu, X., Ma, Y., Sun, W.: Towards real scene super-resolution with raw images. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1723–1731 (2019). <https://doi.org/10.1109/CVPR.2019.00182> [3](#)
88. Xu, X., Ma, Y., Sun, W., Yang, M.H.: Exploiting raw images for real-scene super-resolution. arXiv preprint arXiv:2102.01579 (2021) [3](#)
89. Yin, X., Kolouri, S., Rohde, G.K.: Gat: Generative adversarial training for adversarial example detection and robust classification. In: International Conference on Learning Representations (2019) [2](#)
90. Yu, K., Li, Z., Peng, Y., Loy, C.C., Gu, J.: Reconfigisp: Reconfigurable camera image processing pipeline. ArXiv [abs/2109.04760](#) (2021) [4](#)
91. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: European conference on computer vision. pp. 649–666. Springer (2016) [1](#)
92. Zhang, X., Chen, Q., Ng, R., Koltun, V.: Zoom to learn, learn to zoom. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3762–3770 (2019) [3](#)
93. Zheng, H., Zhang, Z., Gu, J., Lee, H., Prakash, A.: Efficient adversarial training with transferable adversarial examples. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1181–1190 (2020) [5](#)