

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 into making 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 propose 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 Processing

1 Introduction

The most successful methods for a wide range of computer vision tasks rely on deep neural networks [10, 29, 30, 34, 89] (DNNs), including classification, detection, segmentation, scene understanding, scene reconstruction, and generative tasks. Although we rely on the predictions of DNNs for safety-critical applications in robotics, self-driving vehicles, medical diagnostics, and video security, existing networks have been shown to be vulnerable to adversarial attacks [73]: small perturbations to images that are imperceptible to the human vision system can deceive DNNs to make incorrect predictions [51, 55, 62, 72, 77]. As such, defending against adversarial attacks [4, 50, 51, 59, 84] can help resolve failure

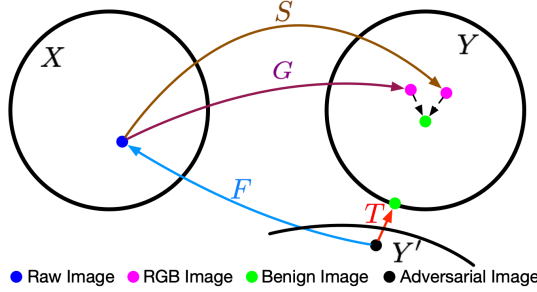


Fig. 1: Existing defense approaches learn an RGB-to-RGB projection from an adversarial distribution (Y') to its natural RGB distribution (Y): $T : Y' \rightarrow Y$. In contrast, our approach learns a mapping via the intermediate natural RAW distribution (X), which is achieved by utilizing three specially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

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 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 [23, 79, 80, 80, 87]. 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 captured in a single 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 [37, 49, 66] or rely on conventional image-processing operations [15, 18, 26, 45]. 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 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 input 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. Code and data to reproduce the findings in this work are available [here](#).

Specifically, we make the following contributions:

- We propose the first adversarial defense method that exploits the natural distribution of RAW domain images.
- The proposed method avoids the generation of adversarial training images and can be used as an off-the-shelf preprocessing module for diverse tasks.
- We analyze 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.

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 [38], each addressing a portion of this image reconstruction problem. In a hardware ISP, these modules are proprietary compute units, and their behaviors are unknown to the user. More importantly, the modules are not differentiable [54, 75]. Two lines of work rely on deep-learning-based approaches to cope with the significant drawback.

The first flavor of methods directly replaces the hardware ISP with a deep-learning-based model to target different application scenarios, such as low-light enhancement [8, 9], super-resolution [85, 86, 90], smartphone camera enhancement [14, 33, 67], and ISP replacement [44]. Nevertheless, the deep-learning-based models used by these works contain a massive number of parameters and are computationally expensive. Hence, their application is limited to offline tasks. In contrast, another thread of works focused on searching for the best hardware ISP hyperparameters for different downstream tasks by using deep-learning-based

approaches. Specifically, Tseng *et al.* [75] proposed differentiable proxy functions to model arbitrary ISP pipelines and used them to find the best hardware ISP hyperparameters for different downstream tasks. Yu *et al.* [88] proposed Re-configISP, which uses different proxy functions for each module of a hardware ISP instead of the entire ISP pipeline. Mosleh *et al.* [54] 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 level of access to target networks, adversarial attacks can be broadly categorized into white-box attacks and black-box attacks.

Among white-box attacks, one important direction is gradient-based attacks [24, 41, 51]. 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 [6, 72]. In the black-box setting, only benign images and their class labels are given, meaning attackers can only query the target model. With the free query, Black-box attacks rely on adversarial transferability to train substitute models [31, 58, 60, 70] or directly estimate the target model gradients [12, 13, 78] to generate adversarial examples. To avoid the transferability 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 [5, 43, 56].

In the physical world, adversarial samples are captured by cameras providing image inputs to target networks. A variety of strategies have been developed to guard the effectiveness of the adversarial patterns in the wild [2, 17, 21, 35]. These methods typically assume that the camera acquisition and subsequent hardware processing do not alter the adversarial patterns. Phan *et al.* [61] have recently realized attacks on 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 [51], input transformation [3, 19], defensive distillation [59], dynamic models [81], loss modifications [57], model ensemble [68] and robust architecture [27]. We next analyze the two representative categories of defense methods.

Adversarial Training (AT): The idea of AT is the following: in each training loop, we augment training data with adversarial examples generated by different attacks. AT is known to “overfit” to the attacks “seen” during training and has been demonstrated to be vastly effective in defending those attacks. However, AT does not generalize well to “unseen” attacks [71]. Furthermore, iteratively generating adversarial images is time-consuming, causing 3-30 times longer than standard training before the model converges [69]. Multiple methods

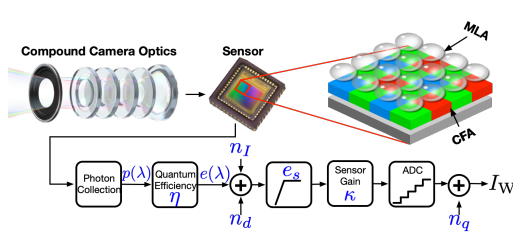


Fig. 2: Overview of the RAW imaging pipeline model. The scene light field is captured by compound camera optics, and then it is gathered by a microlens array layer before being captured on the color filter array. The color-filtered photons are converted into electrons based on quantum efficiency before 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 .

have been proposed to reduce the training time, making AT on large datasets (*e.g.*, ImageNet) possible [23, 79, 80, 91]. Even so, for each specific model, the approach 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. [45].

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 [15, 19, 48], randomization [82], image quilting [26], pixel deflection [63], and deep-learning-based approaches [37, 49, 66]. These IT methods can seamlessly work with different downstream models and tasks. More importantly, 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 existing IT methods to focus on the preprocessing in the RGB distribution, the proposed approach leverages the intermediate natural RAW distribution to remove adversarial patterns.

3 Sensor Image Formation

In this section, we review how a RAW image is formed. When light from the scene enters a camera aperture, it first passes through compound camera optics which focus the light on an image sensor (*e.g.*, CCD and CMOS), where the photons are color-filtered and converted to electrons. Finally, the electrons are converted to digital values, comprising a RAW image. We refer the reader to Karaimer and Brown [39] for a detailed review.

Compound Camera Optics: A compound lens consisting of a sequence of optics is designed to correct optical aberrations. When scene radiance, I_{SCENE} (in the form of a light field) enters a compound lens, the radiance is modulated by the complex optical pipeline and generates the image I_O . Compound optics can be modeled by spatially-varying point spread functions (PSFs) [74].

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 pass 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 to electrons, and the conversion process is specified by the detector quantum efficiency. During this process, electrons can fluctuate randomly which we summarize as electron noise. Two common types of electron noise 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 [74]. 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 quantizes 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 nature 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 exploits 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 samples 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 – image processing pipelines are designed to remove such RAW cues.

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

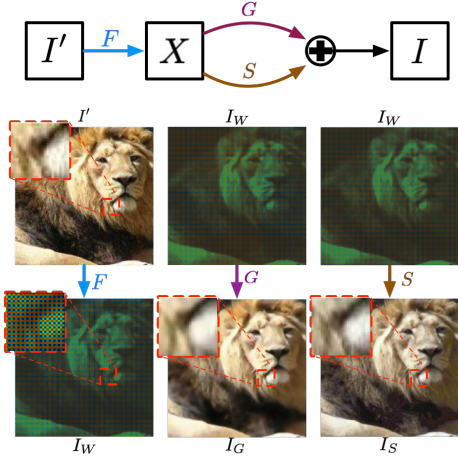


Fig. 3: Overview of the proposed defense approach; see text. We note that the resolution of the RAW image (RGGB) is twice larger than that of the RGB image. We linearly scaled the RAW image in this figure for better visualization.

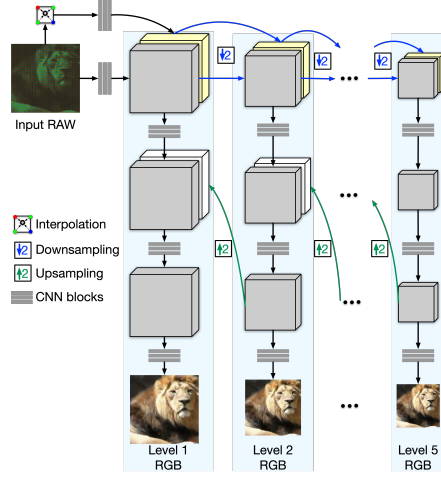


Fig. 4: The architecture of the G operator is adopted and modified from PyNet [32]. The finer operator level exploits upsampled coarser-level features to reconstruct the RGB output. The model is trained sequentially in a coarse-to-fine manner.

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 rely on an additional operator S , a conventional ISP, to our approach, which is implemented as a sequence of cascaded software-based sub-modules. In contrast to the operator F , the operator S is non-differentiable. The 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 are given in the Supplementary Material.

We train this module in a supervised manner with two \mathcal{L}_2 losses. Both \mathcal{L}_2 losses are calculated between the ground truth (GT) RAW and estimated RAW images. The only difference between the two losses is the input RGB image corresponding to the estimated RAW. One is using 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 corresponding RAW distributions. 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 μ and σ to 0 and 1, respectively.

4.2 G Operator: Learned ISP

The G operator is represented by a neural network that converts the I_w generated by the F operator to an RGB image. During this process, we aim to guide local adjustment by global contextual information. This motivates us to devise a pyramidal convolutional neural network to fuse global and local features. To this end, we propose a variant of PyNet [32] consisting of five levels; see Figure 4. Here, the 1st level is the finest, and the 5th level is the coarsest. The finer levels use upsampled features from the coarser levels by concatenating them. We modify PyNet by adding an interpolation layer before the input of each level, interpolating the downsampled RAW Bayer pattern. This facilitates learning as the network only needs to learn the residuals between interpolated RGB and ground truth RGB.

The loss function for this model consists of three components: perceptual, structural similarity, and \mathcal{L}_2 loss. Given an input RAW image, I_W , and the corresponding GT RGB image GT_I , the loss function can be 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 represent the perceptual loss calculated with VGG architecture, structural similarity loss, and \mathcal{L}_2 loss, respectively; β^i and γ^i are the two hyperparameter weights, 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, which is non-differentiable.

While we may use the ISP pipeline of any digital camera we can extract raw and post-ISP data from, we employ a software-based ISP pipeline consisting of the following components: Bayer demosaicing, color balancing, white balancing, contrast improvement, and colorspace conversion sub-modules. Using on the Zurich-Raw-to-RGB dataset [33], we manually tune the hyperparameters of all sub-modules. We refer the reader to the Supplementary Material for a detailed description.

4.4 Operator Training

We use the Zurich-Raw-to-RGB dataset [33] 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 the F and G operators, respectively. We use the following hyperparameters settings: $\omega = 0.7$ in Eq. 2; $\mu = 0$ and $\sigma = 1$ for the Gaussian component 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 $i \in [2, 5]$.

5 Experiments & Analysis

	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 [19]	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 [26]	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 [82]	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 [46]	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 [63]	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 [37]	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 [19]	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 [26]	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 [82]	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 [46]	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 [63]	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 [37]	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.

The proposed method acts as an off-the-shelf input preprocessing module, requiring no additional training to be transferred to different tasks. To validate the effectiveness and generalization capabilities of the proposed defense approach,

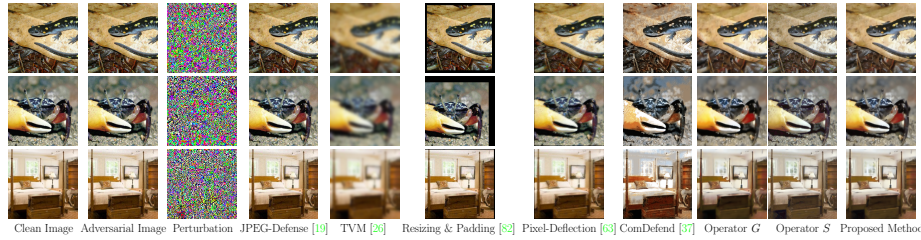


Fig. 5: Qualitative outputs of the proposed method along with both G and S operators and state-of-the-art defense methods on the ImageNet dataset; see text.

	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 [19]	37.41	32.27	24.53	6.21	25.74	10.18	14.12	5.66
TVM [26]	42.64	41.53	45.55	42.24	44.51	38.44	31.88	25.56
HGD [46]	43.39	40.82	44.54	40.88	40.03	39.95	28.61	22.36
Pixel-Deflection [63]	44.13	41.88	46.38	42.32	44.78	37.22	30.73	24.61
ComDefend [37]	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 for mean IoU (mIoU) and mark the best mIoU in bold. Our defense method achieves the best performance in all settings.

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 [25], BIM [42], PGD [52], C&W [7], Newton-Fool [36], and DeepFool [53]. For classification, we use the widely used Foolbox benchmarking suite [64] to implement these attack methods. Since Foolbox does not directly support semantic segmentation and object detection attacks, we use the lightweight TorchAttacks library [40] for generating adversarial examples with FGSM, PGD, and BIM attacks. We also evaluate against the DAG [83] attack, a dedicated attack approach for semantic segmentation and object detection tasks. Moreover, we further evaluate against BPDA [1], an attack method specifically designed for circumventing input transformation defenses that rely on obfuscated gradients. Defending against BPDA with our method, however, requires a slight modification at inference time; see Supplementary Material for

	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 [19]	22.18	17.61	23.40	8.37	20.39	6.55	7.09	4.35
TVM [26]	26.72	24.79	23.91	21.07	20.60	18.53	16.12	14.28
HGD [46]	31.15	26.23	29.44	25.78	24.42	23.08	20.17	19.75
Pixel-Deflection [63]	29.40	25.36	30.28	26.59	28.94	21.30	18.82	17.55
ComDefend [37]	30.62	25.89	28.74	27.69	27.08	21.39	20.83	19.42
Proposed Method	37.83	35.15	39.22	36.20	38.63	35.18	25.09	22.47

Table 3: **Quantitative Comparison to SOTA Input-Transformation Defense Methods on the ADE20K dataset.** We evaluate all methods for mean IoU (mIoU) and mark the best mIoU in bold. Our defense method achieves the best performance in all settings.

	FGSM		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 [19]	39.02	35.88	37.96	33.51	38.85	34.69	30.72	25.07
TVM [26]	48.11	39.66	47.1	44.38	48.94	41.76	39.20	33.18
HGD [46]	50.68	40.06	51.24	45.92	46.80	39.74	41.15	37.23
Pixel-Deflection [63]	53.77	44.82	54.45	47.22	55.32	48.32	46.52	39.87
ComDefend [37]	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 4: **Quantitative Comparison to SOTA Input-Transformation Defenses on the Pascal VOC dataset.** We evaluate all compared methods on mean average precision (mAP) on the Pascal VOC dataset. The best mAPs are marked in bold. Our defense method achieves the best performance in all settings.

details. We note that all applied attacks are untargeted. Definitions of all attack methods are provided in the Supplementary Material.

Baseline Defense Approaches: We compare to the following input transformation defense methods: JPEG compression [19], randomized resizing & padding [82], image quilting [26], TVM [26], HGD [46], pixel deflection [63], and Comdefend [37]. We evaluate all baseline methods on the three vision tasks, except that the randomized resizing & padding 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 [47], ADE20K [92] and Pascal VOC [20] 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 FGSM, 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, for 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 the G and S operators. The G operator learns to mitigate the adversarial pattern, *i.e.*, it recovers a latent image in the presence of severe measurement uncertainty. In contrast, 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: commonly used adversarial attacks and attacks specially designed for attacking semantic segmentation models. For the former, FGSM, BIM, and PGD are used, and we employ DAG [83], a dedicated semantic segmentation attack, for the latter. All attacks are based on pre-trained DeepLabV3 models [11]. 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 and 3. The proposed approach significantly outperforms all baseline methods in all experimental settings. Note that no additional training is required to apply the proposed approach, 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 pre-trained Faster R-CNN [65]. We report the mAP on the Pascal VOC dataset under different experimental settings in Table 4. 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 end-to-end image processing and camera design [16, 76]. 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 5, discussed in the following, 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 6, which we also discuss following, where the defense breaks down from 71% to 53%, when gradually moving from RAW to RGB as intermediate image space. These experiments 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. We next map to other intermediate stages in the processing pipeline, such as demosaicing stage, color balance stage, and 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 6, we observe that the defense performance gradually decreases as we map via a less RAW intermediate space.

Real RAW Versus Synthetic RAW: We further ablate on the dataset used to train our model. Specifically, we trained F and G operators on the Zurich-Raw-to-RGB dataset, HDR-RAW-RGB [28] and MIT-RAW-RGB [22] respectively, and assess how the defense performance changes. Similar to the Zurich-Raw-to-RGB dataset, the RAW images in the HDR-RAW-RGB dataset 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. Experimental results are reported in Table 5. As observed, both RAW distributions of 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.

	FSGM	PGD	C&W	NewtonFool	DeepFool
Zurich-Raw-to-RGB [33]	58.85	66.17	71.40	40.96	72.04
HDR-RAW-RGB [28]	55.57	64.12	71.65	42.36	70.77
MIT-RAW-RGB [22]	40.52	47.29	55.13	28.52	58.49

Table 5: **Quantitative Ablation Study on RAW Training Datasets.** We train F and G operators on 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 6: **Effect of Different Intermediate Mapping Spaces.** We report the Top-1 adversarial defense accuracy on the 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 of exploiting the RAW space in the proposed defense.

5.4 Robustness to Hyperparameter and Operator Deviations

Hyper-parameter ω : We introduce a hyperparameter ω for weighting the contributions of the two operators G and S . Here, we evaluate how varying values of ω affect the overall defense accuracy. As reported in Tab. 7, 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.

Deviations of Operators F and G : The operators F and G are trained separately and used jointly at inference time. We evaluate how deviations in

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 7: **Effect of hyperparameter ω .** We evaluate the impact of the method hyperparameter ω on the effectiveness of the proposed defense method.

	<i>F</i> -300	<i>F</i> -320	<i>F</i> -340	<i>F</i> -360	<i>F</i> -380	<i>F</i> -400
<i>G</i> -300	66.02	66.08	65.93	66.05	66.03	66.05
<i>G</i> -330	66.11	66.04	65.97	66.01	65.99	66.04
<i>G</i> -360	65.98	66.04	66.02	65.89	66.08	69.01
<i>G</i> -390	65.98	65.95	66.00	66.04	65.99	65.94

Table 8: **Robustness to Deviations of F and G .** We evaluate the defense accuracy when mixing operator from different training epochs.

Gaussian Noise σ	0 (no noise)	0.01	0.05	0.1	0.3	0.5
Against FSGM Attack	66.02	66.01	65.98	65.90	65.73	65.64
Against PGD Attack	68.34	68.30	68.22	68.10	68.03	67.83

Table 9: **Robustness to Deviations of F .** We perturb the output of operator F with Gaussian noise of different standard deviations and report the defense accuracy.

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. 8 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. 9 reports that the 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 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 captures – 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 extensively validate the effectiveness of the method on existing classification and segmentation datasets. The 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 diverse vision tasks. We also provide insight into the working principles of the approach, confirming the role of the RAW image space in the proposed method. In the future, we plan to explore RAW natural image statistics as an unsupervised prior for image reconstruction and generative neural rendering tasks.

Acknowledgements: Felix Heide was supported by an NSF CAREER Award (2047359), a Sony Young Faculty Award, Amazon Research Award, and a Project X Innovation Award.

References

1. 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
2. Athalye, A., Engstrom, L., Ilyas, A., Kwok, K.: Synthesizing robust adversarial examples. In: International conference on machine learning. pp. 284–293. PMLR (2018) 4
3. Bahat, Y., Irani, M., Shakhnarovich, G.: Natural and adversarial error detection using invariance to image transformations. arXiv preprint arXiv:1902.00236 (2019) 4
4. 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
5. 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
6. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE Symposium on Security and Privacy (2017) 4
7. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks (2017) 10
8. 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
9. 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
10. 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
11. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation (2017) 12
12. 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
13. 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
14. Dai, L., Liu, X., Li, C., Chen, J.: Awnet: Attentive wavelet network for image isp. In: ECCV Workshops (2020) 3
15. 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

16. 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](#)
17. 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](#)
18. Dziugaite, G.K., Ghahramani, Z., Roy, D.M.: A study of the effect of jpg compression on adversarial images (2016) [2](#)
19. 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](#), [9](#), [10](#), [11](#)
20. 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](#)
21. 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](#)
22. Gharbi, M., Chaurasia, G., Paris, S., Durand, F.: Deep joint demosaicking and denoising. *ACM Transactions on Graphics (TOG)* **35**(6), 191 (2016) [13](#)
23. 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](#)
24. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. *CoRR* **abs/1412.6572** (2015) [4](#)
25. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples (2015) [10](#)
26. Guo, C., Rana, M., Cisse, M., Van Der Maaten, L.: Countering adversarial images using input transformations. *ICLR* (2018) [2](#), [5](#), [9](#), [10](#), [11](#)
27. 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](#)
28. 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](#)
29. 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](#)
30. 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](#)
31. Hu, W., Tan, Y.: Generating adversarial malware examples for black-box attacks based on gan. *ArXiv* **abs/1702.05983** (2017) [4](#)
32. Ignatov, A., Gool, L.V., Timofte, R.: Replacing mobile camera isp with a single deep learning model (2020) [7](#), [8](#)
33. 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](#)
34. 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](#)

35. 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](#)
36. 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](#)
37. 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](#), [9](#), [10](#), [11](#)
38. Karaimer, H.C., Brown, M.S.: A software platform for manipulating the camera imaging pipeline. In: ECCV (2016) [3](#)
39. 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](#)
40. Kim, H.: Torchattacks: A pytorch repository for adversarial attacks (2021) [10](#)
41. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world. arXiv preprint arXiv:1607.02533 (2016) [4](#)
42. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world (2017) [10](#)
43. 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](#)
44. 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](#)
45. 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](#)
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 Conference on Computer Vision and Pattern Recognition. pp. 1778–1787 (2018) [9](#), [10](#), [11](#)
47. 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](#)
48. 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](#)
49. 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](#)
50. 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](#)
51. 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](#)

52. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks (2019) [10](#)
53. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks (2016) [10](#)
54. 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](#)
55. Nakkiran, P.: Adversarial robustness may be at odds with simplicity. arXiv preprint arXiv:1901.00532 (2019) [1](#)
56. 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](#)
57. Pang, T., Xu, K., Dong, Y., Du, C., Chen, N., Zhu, J.: Rethinking softmax cross-entropy loss for adversarial robustness. ICLR (2020) [4](#)
58. 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 CCS 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](#)
59. 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](#)
60. 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](#)
61. 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](#)
62. 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](#)
63. 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](#), [9](#), [10](#), [11](#)
64. Rauber, J., Brendel, W., Bethge, M.: Foolbox: A python toolbox to benchmark the robustness of machine learning models (2018) [10](#)
65. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks (2016) [12](#)
66. Samangouei, P., Kabkab, M., Chellappa, R.: Defense-gan: Protecting classifiers against adversarial attacks using generative models. ICLR (2018) [2](#), [5](#)
67. 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](#)
68. Sen, S., Ravindran, B., Raghunathan, A.: Empir: Ensembles of mixed precision deep networks for increased robustness against adversarial attacks. ICLR (2020) [4](#)

69. 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) [4](#)
70. 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](#)
71. 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) [4](#)
72. 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](#)
73. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks (2014) [1](#)
74. 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) [5](#), [6](#)
75. 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](#)
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. ACM Trans. Graph. **38**(4), 27–1 (2019) [12](#)
77. 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](#)
78. Tu, C.C., Ting, P., Chen, P.Y., Liu, S., Zhang, H., Yi, J., Hsieh, C.J., Cheng, S.M.: Autozoom: 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](#)
79. 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](#)
80. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training. ICLR (2020) [2](#), [5](#)
81. 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](#)
82. Xie, C., Wang, J., Zhang, Z., Ren, Z., Yuille, A.: Mitigating adversarial effects through randomization. arXiv preprint arXiv:1711.01991 (2017) [5](#), [9](#), [10](#), [11](#)
83. Xie, C., Wang, J., Zhang, Z., Zhou, Y., Xie, L., Yuille, A.: Adversarial examples for semantic segmentation and object detection (2017) [10](#), [12](#)
84. 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](#)

85. 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**
86. Xu, X., Ma, Y., Sun, W., Yang, M.H.: Exploiting raw images for real-scene super-resolution. arXiv preprint arXiv:2102.01579 (2021) **3**
87. 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**
88. Yu, K., Li, Z., Peng, Y., Loy, C.C., Gu, J.: Reconfigisp: Reconfigurable camera image processing pipeline. ArXiv **abs/2109.04760** (2021) **4**
89. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: European conference on computer vision. pp. 649–666. Springer (2016) **1**
90. 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**
91. 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**
92. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 633–641 (2017) **11**