

# What Else Can Fool Deep Learning?

## Addressing Color Constancy Errors on Deep Neural Network Performance

Mahmoud Afifi<sup>1</sup>  
<sup>1</sup>York University, Toronto  
mafifi@eecs.yorku.ca

Michael S Brown<sup>1,2</sup>  
<sup>2</sup>Samsung AI Center, Toronto  
mbrown@eecs.yorku.ca

### Abstract

There is active research targeting local image manipulations that can fool deep neural networks (DNNs) into producing incorrect results. This paper examines a type of global image manipulation that can produce similar adverse effects. Specifically, we explore how strong color casts caused by incorrectly applied computational color constancy – referred to as white balance (WB) in photography – negatively impact the performance of DNNs targeting image segmentation and classification. In addition, we discuss how existing image augmentation methods used to improve the robustness of DNNs are not well suited for modeling WB errors. To address this problem, a novel augmentation method is proposed that can emulate accurate color constancy degradation. We also explore pre-processing training and testing images with a recent WB correction algorithm to reduce the effects of incorrectly white-balanced images. We examine both augmentation and pre-processing strategies on different datasets and demonstrate notable improvements on the CIFAR-10, CIFAR-100, and ADE20K datasets.

### 1. Introduction

There is active interest in local image manipulations that can be used to fool deep neural networks (DNNs) into producing erroneous results. Such “adversarial attacks” often result in drastic misclassifications. We examine a less explored problem of *global* image manipulations that can result in similar adverse effects on DNNs’ performance. In particular, we are interested in the role of computational color constancy, which makes up the white-balance (WB) routine on digital cameras.

We focus on computational color constancy because it represents a common source of global image errors found in real images. When WB is applied incorrectly on a camera, it results in an undesirable color cast in the captured image. Images with such strong color casts are often discarded

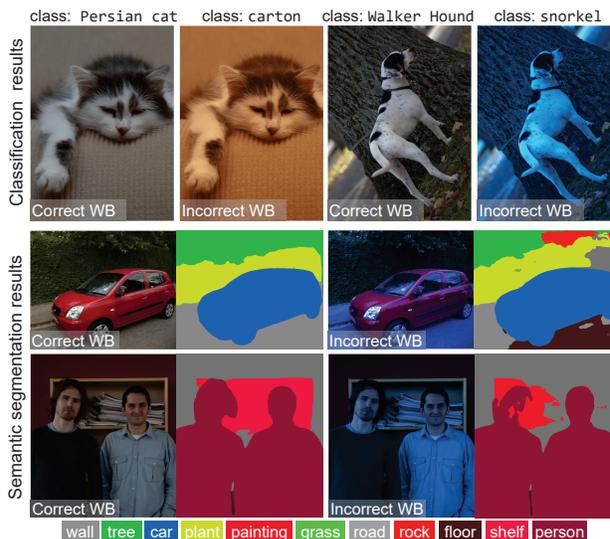


Figure 1. The effect of correct/incorrect computational color constancy (i.e., white balance) on (top) classification results by ResNet [29]; and (bottom) semantic segmentation by RefineNet [39].

by users. As a result, online image databases and repositories are biased to contain mostly correctly white-balanced images. This is an implicit assumption that is not acknowledged for datasets composed of images crawled from the web and online. However, in real-world applications, it is unavoidable that images will, at some point, be captured with the incorrect WB applied. Images with incorrect WB can have unpredictable results on DNNs trained on white-balanced biased training images, as demonstrated in Fig. 1.

**Contribution** We examine how errors related to computational color constancy can adversely affect DNNs focused on image classification and semantic segmentation. In addition, we show that image augmentation strategies used to expand the variation of training images are not well suited to mimic the type of image degradation caused by color constancy errors. To address these problems, we introduce a novel augmentation method that can accurately emulate

realistic color constancy degradation. We also examine a newly proposed WB correction method [2] to pre-process testing and training images. Experiments on CIFAR-10, CIFAR-100, and the ADE20K datasets using the proposed augmentation and pre-processing correction demonstrate notable improvements to test image inputs with color constancy errors.

## 2. Related Work

**Computational Color Constancy** Cameras have on-board image signal processors (ISPs) that convert the raw-RGB sensor values to a standard RGB output image (denoted as an sRGB image) [33, 47]. Computational color constancy, often referred to as WB in photography, is applied to mimic the human’s ability to perceive objects as the same color under any type of illumination. WB is used to identify the color temperature of the scene’s illumination either manually or automatically by estimating the scene’s illumination from an input image (e.g., [1, 6, 7, 9, 17, 25, 30, 51]). After WB is applied to the raw-RGB image, a number of additional nonlinear photo-finishing color manipulations are further applied by the ISP to *render* the final sRGB image [2]. These photo-finishing operations include, but are not limited to, hue/saturation manipulation, general color manipulation, and local/global tone mapping [8, 27, 33, 44, 47]. Cameras generally have multiple photo-finishing styles the user can select [2, 33, 34].

**Post-WB Correction in sRGB Images** When WB is applied incorrectly, it results in sRGB images with strong color casts. Because of the nonlinear photo-finishing operations applied by the ISP after WB, correcting mistakes in the sRGB image is non-trivial [2, 45]. Current solutions require meta-data, estimated from radiometric calibration or raw-image reconstruction methods (e.g., [14, 34, 45]), that contains the necessary information to undo the particular nonlinear photo-finishing processes applied by the ISP. By converting back to a raw-RGB space, the correct WB can be applied using a diagonal correction matrix and then re-rendered by the ISP. Unfortunately, meta-data to inverse the camera pipeline and re-render the image is rarely available, especially for sRGB images gleaned from the web—as is the case with existing computer vision datasets. Recently, it was shown that white balancing sRGB images can be achieved by estimating a high-degree polynomial correction matrix [2]. The work in [2], referred to WB for sRGB images (WB-sRGB), introduces a data-driven framework to estimate such polynomial matrix for a given testing image. We build on the WB-sRGB [2] by extending this framework to emulate WB errors on the final sRGB images, instead of correcting WB. We also used the WB-sRGB method [2] to examine applying a pre-process WB correction on training

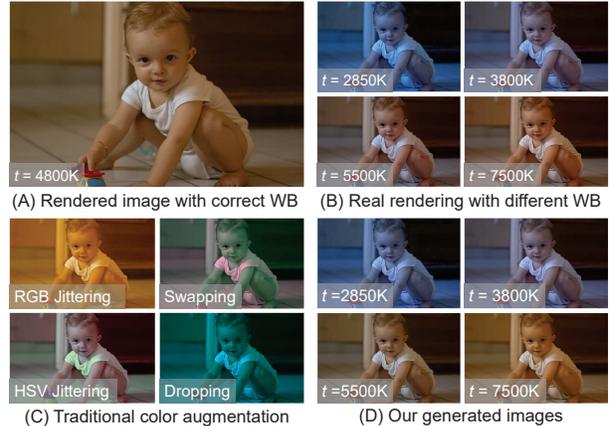


Figure 2. (A) An sRGB image from a camera with the correct WB applied. (B) Images from the same camera with the incorrect WB color temperatures ( $t$ ) applied. (C) Images generated by processing image (A) using existing augmentation methods—the images clearly do not represent those in (B). (D) Images generated from (A) using our proposed method detailed in Sec. 4.

and testing images in order to improve the performance of DNN models against incorrectly white-balanced images.

**Adversarial Attacks** DNN models are susceptible to adversarial attacks in the form of *local* image manipulation (e.g., see [18, 26, 37, 54]). These images are created by adding a carefully crafted imperceptible perturbation layer to the original image [26, 54]. Such perturbation layers are usually represented by *local* non-random adversarial noise [3, 26, 41, 54, 58] or *local* spatial transformations [57]. Adversarial examples are able to misguide pre-trained models to predict either a certain wrong response (i.e., targeted attack) or any wrong response (i.e., untargeted attack) [3, 12, 40]. While incorrect color constancy is not an explicit attempt at an adversarial attack, the types of failures produced by this *global* modification act much like an untargeted attack and can adversely affect DNNs’ performance.

**Data Augmentation** To overcome limited training data and to increase the visual variation, image augmentation techniques are applied to training images. Existing image augmentation techniques include: geometric transformations (e.g., rotation, translation, shearing) [19, 28, 28, 46], synthetic occlusions [60], pixel intensity processing (e.g., equalization, contrast adjustment, brightness, noise) [19, 56], and color processing (e.g., RGB color jittering and PCA-based shifting, HSV jittering, color channel dropping, color channel swapping) [15, 19, 23, 32, 36, 38, 42, 48, 49]. Traditional color augmentation techniques randomly change the original colors of training images aiming for better generalization and robustness of the trained model in

the inference phase. However, existing color augmentation methods often generate unrealistic colors which rarely happen in reality (e.g., green skin or purple grass). More importantly, the visual appearance of existing color augmentation techniques does not well represent the color casts produced by incorrect WB applied onboard cameras, as shown in Fig. 2. As demonstrated in [4, 13, 22], image formation has an important effect on the accuracy of different computer vision tasks. Recently, a simplified version of the camera imaging pipeline was used for data augmentation [13]. This augmentation method in [13], however, explicitly did not consider the effects of incorrect WB due to the subsequent nonlinear operations applied after WB. To address this issue, we propose a camera-based augmentation technique that can synthetically generate images with realistic WB settings.

**DNN Normalization Layers** Normalization layers are commonly used to improve the efficiency of the training process. Such layers apply simple statistics-based shifting and scaling operations to the activations of network layers. The shift and scale factors can be computed either from the entire mini-batch (i.e., batch normalization [31]) or from each training instance (i.e., instance normalization [55]). Recently, batch-instance normalization (BIN) [43] was introduced to ameliorate problems related to styles/textures in training images by balancing between batch and instance normalizations based on the current task. Though the BIN is designed to learn the trade-off between keeping or reducing original training style variations using simple statistics-based operations, the work in [43] does not provide any study regarding incorrect WB settings. The augmentation and pre-processing methods proposed in our work directly target training and testing images and do not require any change to a DNNs architecture or training regime.

### 3. Effects of WB Errors on Pre-trained DNNs

We begin by studying the effect of incorrectly white-balanced images on pre-trained DNN models for image classification and semantic segmentation. As a motivation, Fig. 3 shows two different WB settings applied to the same image. Fig. 3 shows that the DNN’s attention for the same scene is considerably altered by changing the WB setting.

For quantitative evaluations, we adopted several DNN models trained for the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2012 [21] and the ADE20K Scene Parsing Challenge 2016 [61]. Generating an entirely new labeled testing set composed of images with incorrect WB is an enormous task—ImageNet classification includes 1,000 classes and pixel-accurate semantic annotation requires ~60 minutes per image [50]. In lieu of a new testing set, we applied our method which emulates WB errors to

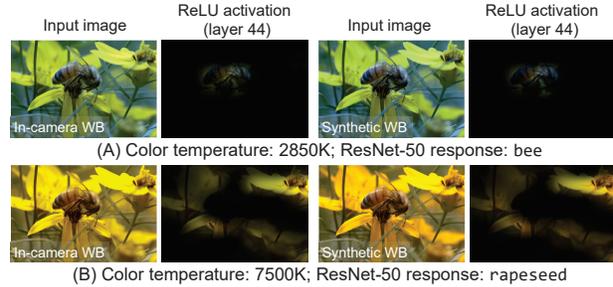


Figure 3. Image rendered with two different color temperatures (denoted by  $t$ ) using in-camera rendering and our method. (A) Image class is `bee`. (B) Image class is `rapeseed`. Classification results were obtained by ResNet-50 [29].

the validation images of each dataset. Our method will be detailed shortly in Sec. 4.

**Classification** We apply our method to ImageNet’s validation set to generate images with five different color temperatures and two different photo-finishing styles for a total of ten WB variations for each validation image; 899 grayscale images were excluded from this process. In total, we generated 491,010 images. We examined the following six well-known DNN models, trained on the original ImageNet training images: AlexNet [36], VGG-16 & VGG-19 [52], GoogLeNet [53], and ResNet-50 & ResNet-101 [29]. Table 1 shows the accuracy drop for each model when tested on our generated validation set (i.e., with different WB and photo-finishing settings) compared to the original validation set. In most cases, there is a drop of ~10% in accuracy. Fig. 4 shows an example of the impact of incorrect WB.

**Semantic Segmentation** We used the ADE20K validation set for 2,000 images, and generated ten images with different WB/photo-finishing settings for each image. At the end, we generated a total of 20,000 new images. We tested the following two DNN models trained on the original ADE20K training set: DilatedNet [16, 59] and RefineNet [39]. Table 2 shows the effect of improperly white-balanced images on the intersection-over-union (IoU) and



Figure 4. Pre-trained models are negatively impacted by incorrect WB settings. (A) Original image. (B) Generated images with different WB color temperatures (denoted by  $t$ ). Classification results of: VGG-16 [52], GoogLeNet [53], and ResNet-50 [29] are written on top of each image. The terms E and S stand for Egyptian and Siamese, respectively.

Table 1. Adverse performance on ImageNet [21] due to the inclusion of incorrect WB versions of its validation images. The models were trained on the original ImageNet training set. The reported numbers denote the changes in the top-1 accuracy achieved by each model.

Model	Effect on top-1 accuracy
AlexNet [36]	-0.112
VGG-16 [52]	-0.104
VGG-19 [52]	-0.102
GoogLeNet [53]	-0.107
ResNet-50 [29]	-0.111
ResNet-101 [29]	-0.109

Table 2. Adverse performance on ADE20K [61] due to the inclusion of incorrect WB versions of its validation images. The models were trained on ADE20K’s original training set. The reported numbers denote the changes in intersection-over-union (IoU) and pixel-wise accuracy (pxl-acc) achieved by each model on the original validation.

Model	Effect on IoU	Effect on pxl-acc
DilatedNet [16, 59]	-0.023	-0.024
RefineNet [39]	-0.031	-0.026

pixel-wise accuracy (pxl-acc) obtained by the same models on the original validation set. While DNNs for segmentation fare better than the results for classification, we still incur a drop of over 2% in performance.

## 4. Proposed Method to Emulate WB Errors

Given an sRGB image, denoted as  $\mathbf{I}_{t_{\text{corr}}}$ , that is assumed to be white-balanced with the correct color temperature, our goal is to modify  $\mathbf{I}_{t_{\text{corr}}}$ ’s colors to mimic its appearance as if it were rendered by a camera with different (incorrect) color temperatures,  $t$ , under different photo-finishing styles. Since we do not have access to  $\mathbf{I}_{t_{\text{corr}}}$ ’s original raw-RGB image, we cannot re-render the image from raw-RGB to sRGB using a standard camera pipeline. Instead, we have adopted a data-driven method that mimics this manipulation directly in the sRGB color space. Our framework draws heavily from the WB-sRGB data-driven framework [2], which was proposed to correct improperly white-balanced sRGB images. Our framework, however, “emulates” WB errors on the rendered sRGB images. Fig. 5 provides an overview of our method.

### 4.1. Dataset

Our method relies on a large dataset of sRGB images generated by [2]. This dataset contains images rendered with different WB settings and photo-finishing styles. There is a ground truth sRGB image (i.e., rendered with the “correct” color temperature) associated with each training image. The training sRGB images were rendered using five different color temperatures: 2850 Kelvin (K), 3800K, 5500K, 6500K, and 7500K. In addition, each image was rendered using different camera photo-finishing styles. In our WB emulation framework, we used 17,970 images from

this dataset (1,797 correct sRGB images each with ten corresponding images rendered with five different color temperatures and two different photo-finishing styles, Camera Standard and Adobe Standard).

### 4.2. Color Mapping

Next, we compute a mapping between the correct white-balanced sRGB image to each of its ten corresponding images. We follow the same procedure of the WB-sRGB method [2] and use a kernel function,  $\varphi$ , to project RGB colors into a high-dimensional space. Then, we perform polynomial data fitting on these projected values. Specifically, we used  $\varphi: [\mathbf{R}, \mathbf{G}, \mathbf{B}]^T \rightarrow [\mathbf{R}, \mathbf{G}, \mathbf{B}, \mathbf{R}\mathbf{G}, \mathbf{R}\mathbf{B}, \mathbf{G}\mathbf{B}, \mathbf{R}^2, \mathbf{G}^2, \mathbf{B}^2]^T$  [24]. The data fitting can be represented by a color transformation matrix  $\mathbf{M}_{t_{\text{corr}} \rightarrow t}$  computed by the following minimization equation:

$$\arg \min_{\mathbf{M}_{t_{\text{corr}} \rightarrow t}} \|\mathbf{M}_{t_{\text{corr}} \rightarrow t} \varphi(\mathbf{I}_{t_{\text{corr}}}) - \mathbf{I}_t\|_F, \quad (1)$$

where  $\mathbf{I}_{t_{\text{corr}}}$  and  $\mathbf{I}_t$  are  $3 \times n$  color matrices of the white-balanced image rendered with the correct color temperature  $t_{\text{corr}}$  and color values of the same image rendered with the target different color temperature  $t$ , respectively,  $n$  is the total number of pixels in each image,  $\|\cdot\|_F$  is the Frobenius norm, and  $\mathbf{M}_{t_{\text{corr}} \rightarrow t}$  is represented as a nonlinear  $3 \times 9$  full matrix.

We compute a color transformation matrix between each pair of correctly white-balanced image and its corresponding target image rendered with a specific color temperature and photo-finishing. In the end, we have *ten* matrices associated with each image in our training data.

### 4.3. Color Feature

As shown in Fig. 5, when augmenting an input sRGB image to have different WB settings, we search our dataset for similar sRGB images to the input image. This search is not based on scene content, but on the color distribution of the image. As a result, we represent each image in the training set with the RGB-*uv* projected color histogram feature used in [2]. Each histogram feature is represented as an  $m \times m \times 3$  tensor. To further reduce the size of the histogram feature, we apply principal component analysis (PCA) to the three-layer histogram feature. This transformation maps the zero-centered vectorized histogram to a new lower-dimensional space. Our implementation used a 55-dimensional PCA vector. Our final training data therefore consists of the compacted feature vector of each training white-balanced image, the associated color transformation matrices, and the PCA coefficient matrix and bias vector.

### 4.4. KNN Retrieval

Given a new input image  $\mathbf{I}_{\text{in}}$ , we extract its compacted color feature  $\mathbf{v}$ , and then search for training examples with

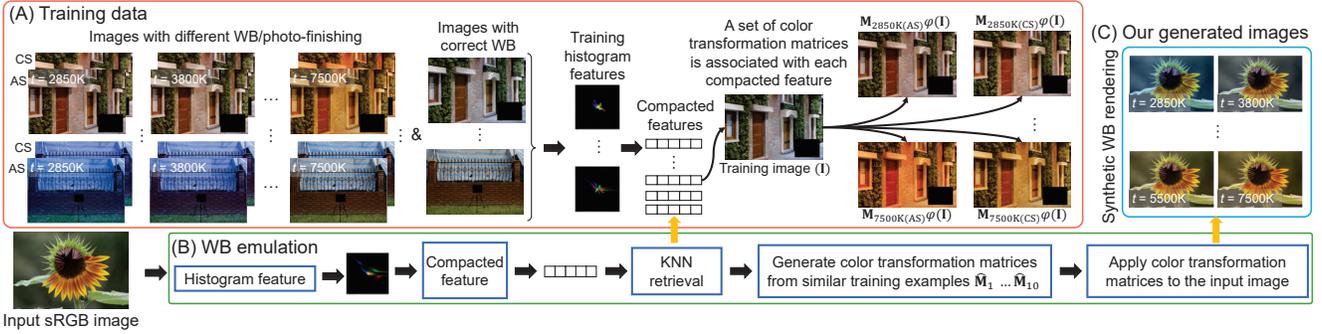


Figure 5. Our WB emulation framework. (A) A dataset of 1,797 correctly white-balanced sRGB images [2]; each image has ten corresponding sRGB images rendered with five different color temperatures and two photo-finishing styles, Camera Standard (CS) and Adobe Standard (AS). For each white-balanced image, we generate its compact histogram feature and ten color transformation matrices to the corresponding ten images. (B) Our WB emulation pipeline (detailed in Sec. 4). (C) The augmented images for the input image that represent different color temperatures (denoted by  $t$ ) and photo-finishing styles.

color distributions similar to the input image’s color distribution. The  $L_2$  distance is adopted as a similarity metric between  $\mathbf{v}$  and the training compacted color features. Afterwards, we retrieve the color transformation matrices associated with the nearest  $k$  training images. The retrieved set of matrices is represented by  $\mathbf{M}_s = \{\mathbf{M}_s^{(j)}\}_{j=1}^k$ , where  $\mathbf{M}_s^{(j)}$  represents the color transformation matrix that maps the  $j^{\text{th}}$  white-balanced training image colors to their corresponding image colors rendered with color temperature  $t$ .

#### 4.5. Transformation Matrix

After computing the distance vector  $\mathbf{d}$  between  $\mathbf{v}$  and the nearest training features, we compute a weighting vector  $\alpha$  as follows [2]:

$$\alpha_j = \frac{\exp(-\mathbf{d}_j^2/2\sigma^2)}{\sum_{k'=1}^k \exp(-\mathbf{d}_{k'}^2/2\sigma^2)}, j \in [1, \dots, k], \quad (2)$$

where  $\sigma$  is the radial basis function parameter. We used  $\sigma = 0.25$  in our experiments. We construct the final color transformation matrix  $\hat{\mathbf{M}}_{t_{\text{corr}} \rightarrow t}$  as a linear weighted combination of the retrieved color transformation matrices  $\mathbf{M}_s$ . This process is performed as follows [2]:

$$\hat{\mathbf{M}}_{t_{\text{corr}} \rightarrow t} = \sum_{j=1}^k \alpha_j \mathbf{M}_s^{(j)}. \quad (3)$$

Lastly, the “re-rendered” image  $\hat{\mathbf{I}}_t$  with color temperature  $t$  is computed as:

$$\hat{\mathbf{I}}_t = \hat{\mathbf{M}}_{t_{\text{corr}} \rightarrow t} \varphi(\mathbf{I}_{\text{in}}). \quad (4)$$

### 5. Experiments

**Robustness Strategies** Our goal is to improve the performance of DNN methods in the face of test images that

may have strong global color casts due to computational color constancy errors. Based on the WB-sRGB framework [2] and the modified framework discussed in Sec. 4, we examine three strategies to improve the robustness of the DNN models.

(1) The first strategy is to apply a WB correction to each testing image in order to remove any unexpected color casts during the inference time. Note that this approach implicitly assumes that the training images are correctly WB. In our experiments, we used the WB-sRGB method [2] to correct the test images, because it currently achieves the state-of-the-art on white balancing sRGB rendered images. We examined adapting the simple diagonal-based correction – which is applied by traditional WB methods that are intended to be applied on raw-RGB images (e.g., gray-world [10]) – but found that they give inadequate results when applied on sRGB images, as also demonstrated in [2]. In fact, applying diagonal-based correction directly on the training image is similar to multiplicative color jittering. This is why we need to use a nonlinear color manipulation (e.g., polynomial correction estimated by [2]) for more accurate WB correction for sRGB images. An example of the difference is shown in Fig. 6.

It is worth mentioning that the training data used by the WB-sRGB method has five fixed color temperatures (2850K, 3800K, 5500K, 6500K, 7500K), all with color correction matrices mapping to their corresponding correct WB. In most cases, one of these five fixed color temperatures will be visually similar to the correct WB. Thus, if the WB-sRGB method is applied to an input image that is already correctly white-balanced, the computed transformation will act as an identity.

(2) The second strategy considers the case that some of the training images may include some incorrectly

white-balanced images. We, therefore, also apply the WB correction step to all the training images as well as testing images. This again uses the WB-sRGB method [2] on both testing and training images.

(3) The final strategy is to augment the training dataset based on our method described in Sec. 4. Like other augmentation approaches, there is no pre-processing correction required. The assumption behind this augmentation process is that the robustness of DNN models can be improved by training on augmented images that serve as exemplars for color constancy errors.

**Testing Data Categories** Testing images are grouped into two categories. In Category 1 (Cat-1), we expand the original testing images in the CIFAR-10, CIFAR-100, and ADE20K datasets by applying our method to emulate camera WB errors (described in Sec. 4). Each test image now has ten (10) variations that share the same ground truth labels. We acknowledge this is less than optimal, given that the same method to modify the testing image is used to augment the training images. However, we are confident in the proposed method’s ability to emulate WB errors that we feel Cat-1 images represents real-world examples. With that said, we do not apply strategies 1 and 2 to Cat-1, as the WB-sRGB method is based on a similar framework used to generate the testing images. For the sake of completeness, we also include Category 2 (Cat-2), which consists of new datasets generated directly from raw-RGB images. Specifically, raw-RGB images are rendered using the full in-camera pipeline to sRGB images with in-camera color constancy errors. As a result, Cat-2’s testing images exhibit accurate color constancy errors but contain fewer testing images for which we have provided the ground truth labels.



Figure 6. (A) Images with different categories of “dogs” rendered with incorrect WB settings. (B) Corrected images using gray-world (GW) [10]. (C) Corrected images using the WB-sRGB method [2]. Predicted class by AlexNet is written on top of each image. Images in (A) and (B) are misclassified.

## 5.1. Experimental Setup

We compare the three above strategies with two existing and widely adopted color augmentation processes: RGB color jittering and HSV jittering.

**Our Method** The nearest neighbor searching was applied using  $k = 25$ . The proposed WB augmentation model runs in 7.3 sec (CPU) and 1.0 sec (GPU) to generate ten 12-mega-pixel images. The reported runtime was computed using Intel® Xeon® E5-1607 @ 3.10 GHz CPU and NVIDIA™ Titan X GPU.

**Existing Color Augmentation** To the best of our knowledge, there is no standardized approach for existing color augmentation methods. Accordingly, we tested different settings and selected the settings that produce the best results.

For RGB color jittering, we generated ten images with new colors by applying a random shift  $x \sim \mathcal{N}(\mu_x, \sigma^2)$  to each color channel of the image. For HSV jittering, we generated ten images with new colors by applying a random shift  $x$  to the hue channel and multiplying each of the saturation and value channels by a random scaling factor  $s \sim \mathcal{N}(\mu_s, \sigma^2)$ . We found that  $\mu_x = -0.3$ ,  $\mu_s = 0.7$ , and  $\sigma = 0.6$  give us the best compromise between having color diversity with low color artifacts during the augmentation process.

## 5.2. Network Training

For image classification, training new models on ImageNet dataset requires unaffordable efforts—for instance, ILSVRC 2012 consists of  $\sim 1$  million images and would be  $\sim 10$  million images after applying any of the color augmentation techniques. For that reason, we perform experiments on CIFAR-10 and CIFAR-100 datasets [35] due to a more manageable number of images in each dataset.

We trained SmallNet [46] from scratch on CIFAR-10. We also fine-tuned AlexNet [36] to recognize the new classes in CIFAR-10 and CIFAR-100 datasets. For semantic segmentation, we fine-tuned SegNet [5] on the training set of the ADE20K dataset [61].

We train each model on: (i) the original training images, (ii) the WB-sRGB method [2] applied to the original training images, and (iii) original training images with the additional images produced by color augmentation methods. For color augmentation, we examined RGB color jittering, HSV jittering, and our WB augmentation. Thus, we trained five models for each CNN architecture, each of which was trained on one of the mentioned training settings.

For fair comparisons, we trained each model for the same number of iterations. Specifically, the training was for  $\sim 29,000$  and  $\sim 550,000$  iterations for image classification

and semantic segmentation tasks, respectively. We adjusted the number of epochs to make sure that each model was trained on the same number of mini-batches for fair comparison between training on augmented and original sets. Note that by using a fixed number of iterations to train models with both original training data and augmented data, we did not fully exploit the full potential of the additional training images when we trained models using additional augmented data.

The training was performed using NVIDIA™ Titan X GPU. The details of training parameters are given in supplemental materials.

### 5.3. Results on Cat-1

Cat-1 tests each model using test images that have been generated by our method described in Sec. 4.

**Classification** We used the CIFAR-10 testing set (10,000 images) to test SmallNet and AlexNet models trained on the training set of the same dataset. We also used the CIFAR-100 testing set (10,000 images) to evaluate the AlexNet model trained on CIFAR-100. After applying our WB emulation to the testing sets, we have 100,000 images for each testing set of CIFAR-10 and CIFAR-100. The top-1 accuracies obtained by each trained model are shown in Table 3. The best results on our expanded testing images, which include strong color casts, were obtained using models trained on our proposed WB augmented data.

Interestingly, the experiments show that applying WB correction [2] on the training data, in most cases, improves the accuracy using both the original and expanded test sets. DNNs that were trained on WB augmented training images achieve the best improvement on the original testing images compared to using other color augmenters.

**Semantic Segmentation** We used the ADE20K validation set using the same setup explained in Sec. 3. Table 4 shows the obtained pxl-acc and IoU of the trained SegNet models. The best results were obtained with our WB augmentation; Fig. 7 shows qualitative examples. Additional examples are also given in supplemental materials.

### 5.4. Results on Cat-2

Cat-2 data requires us to generate and label our own testing image dataset using raw-RGB images. To this end, we collected 518 raw-RGB images containing CIFAR-10 object classes from the following datasets: HDR+ Burst Photography dataset [27], MIT-Adobe FiveK dataset [11], and Raise dataset [20]. We rendered all raw-RGB images with different color temperatures and two photo-finishing styles using the Adobe Camera Raw module. Adobe Camera Raw accurately emulates the ISP onboard a camera

Table 3. [Cat-1] Results of SmallNet [46] and AlexNet [36] on CIFAR dataset [35]. The shown accuracies obtained by models trained on: original training, “white-balanced”, and color augmented sets. The testing was performed using: original testing set and testing set with different synthetic WB settings (denoted as diff. WB). The results of the baseline models (i.e., trained on the original training set) are highlighted in green, while the best result for each testing set is shown bold. We highlight best results obtained by color augmentation techniques in yellow. Effects on baseline model results are shown in parentheses.

Cat-1		SmallNet [46] on CIFAR-10 [35]	
Training set	Original	Diff. WB	
Original training set	0.799	0.655	
“White-balanced” set	0.801 (+0.002)	0.683 (+0.028)	
HSV augmented set	0.801 (+0.002)	0.747 (+0.092)	
RGB augmented set	0.780 (-0.019)	0.765 (+0.11)	
WB augmented set (ours)	<b>0.809 (+0.010)</b>	<b>0.786 (+0.131)</b>	
Cat-1		AlexNet [36] on CIFAR-10 [35]	
Original training set	0.933	0.797	
“White-balanced” set	0.932 (-0.001)	0.811 (+0.014)	
HSV augmented set	0.923 (-0.010)	0.864 (+0.067)	
RGB augmented set	0.922 (-0.011)	0.872 (+0.075)	
WB augmented set (ours)	0.926 (-0.007)	<b>0.889 (+0.092)</b>	
Cat-1		AlexNet [36] on CIFAR-100 [35]	
Original training set	0.768	0.526	
“White-balanced” set	0.757 (-0.011)	0.543 (+0.017)	
HSV augmented set	0.722 (-0.044)	0.613 (+0.087)	
RGB augmented set	0.723 (-0.045)	0.645 (+0.119)	
WB augmented set (ours)	0.735 (-0.033)	<b>0.670 (+0.144)</b>	

Table 4. [Cat-1] Results of SegNet [5] on the ADE20K validation set [61]. The shown intersection-over-union (IoU) and pixel-wise accuracy (pxl-acc) were achieved by models trained using: original training, “white-balanced”, and color augmented sets. The testing was performed using: original testing set and testing set with different synthetic WB settings (denoted as diff. WB). Effects on results of SegNet trained on the original training set are shown in parentheses. Highlight marks are as described in Table 3.

Cat-1		IoU	
	Original	Diff. WB	
Original training set	0.208	0.180	
“White-balanced” set	<b>0.210 (+0.002)</b>	0.197 (+0.017)	
HSV augmented set	0.192 (-0.016)	0.185 (+0.005)	
RGB augmented set	0.195 (-0.013)	0.190 (+0.010)	
WB augmented set (ours)	0.202 (-0.006)	<b>0.199 (+0.019)</b>	
Cat-1		pxl-acc	
Original training set	0.603	0.557	
“White-balanced” set	<b>0.605 (+0.002)</b>	0.579 (+0.022)	
HSV augmented set	0.583 (-0.020)	0.536 (-0.021)	
RGB augmented set	0.544 (-0.059)	0.534 (-0.023)	
WB augmented set (ours)	0.597 (-0.006)	<b>0.581 (+0.024)</b>	

and produces results virtually identical to what the in-camera processing would produce [2]. Images that contain multiple objects were manually cropped to include only the interesting objects—namely, the CIFAR-10 classes. At the end, we generated 15,098 rendered testing images that reflect real in-camera WB settings. We used the following testing sets in our experiments:

- (i) **In-camera auto WB** contains images rendered with the auto WB (AWB) correction setting in Adobe Camera Raw, which mimics the camera’s AWB functionality.

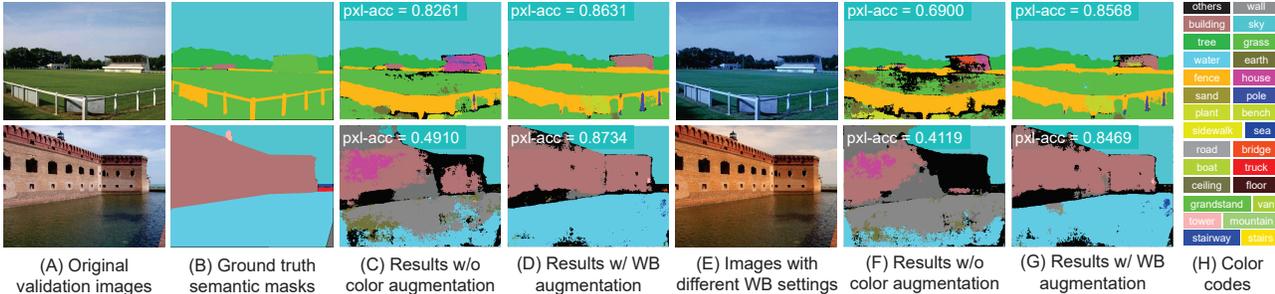


Figure 7. Results of SegNet [5] on the ADE20K validation set [61]. (A) Original validation image. (B) Ground truth semantic mask. (C) & (D) Results of trained model w/o/w color augmentation using image in (A), respectively. (E) Image with a different WB. (F) & (G) Results w/o and with color augmentation using image in (E), respectively. (H) Color codes. The term ‘pxl-acc’ refers to pixel-wise accuracy.

Table 5. [Cat-2] Results of SmallNet [46] and AlexNet [36]. The shown accuracies were obtained using trained models on the original training, “white-balanced”, and color augmented sets. Effects on results of models trained on the original training set are shown in parentheses. Highlight marks are as described in Table 3.

Cat-2 Training Set	SmallNet		
	In-cam AWB	In-cam Diff. WB	WB pre-processing
Original training set	0.467	0.404	0.461
“White-balanced” set	<b>0.496 (+0.029)</b>	0.471 (+0.067)	<b>0.492 (+0.031)</b>
HSV augmented set	0.477 (+0.001)	0.462 (+0.058)	0.481 (+0.02)
RGB augmented set	0.474 (+0.007)	0.475 (+0.071)	0.470 (+0.009)
WB augmented set (ours)	0.494 (+0.027)	<b>0.496 (+0.092)</b>	0.484 (+0.023)
Cat-2	AlexNet		
	In-cam AWB	In-cam Diff. WB	WB pre-processing
Original training set	0.792	0.734	0.772
“White-balanced” set	0.784 (-0.008)	0.757 (+0.023)	0.784 (+0.012)
HSV augmented set	0.790 (+0.002)	0.771 (+0.037)	0.779 (+0.007)
RGB augmented set	0.791 (-0.001)	0.779 (+0.045)	0.783 (+0.011)
WB augmented set (ours)	<b>0.799 (+0.007)</b>	<b>0.788 (+0.054)</b>	<b>0.787 (+0.015)</b>

AWB does fail from time to time; we manually removed images that had a noticeable color cast. This set of images is intended to be equivalent to testing images on existing image classification datasets.

(ii) **In-camera WB settings** contains images rendered with the different color temperatures and photo-finishing styles. This set represents testing images that contain WB color cast errors.

(iii) **WB pre-processing correction applied to set (ii)** contains images of set (ii) after applying the WB-sRGB correction [2]. This set is used to study the potential improvement of applying a pre-processing WB correction in the inference phase.

Table 5 shows the top-1 accuracies obtained by SmallNet and AlexNet on the external testing sets. The experiments show the accuracy is reduced by  $\sim 6\%$  when the testing set is images that have been modified with incorrect WB settings compared with their original accuracies obtained with “properly” white-balanced images using the in-camera AWB. We also notice that the best accuracies are obtained by applying either a pre-processing WB on both training/testing images or our WB augmentation in an end-to-end manner. Examples of misclassified images are shown in Fig. 8. Additional examples are also given in sup-

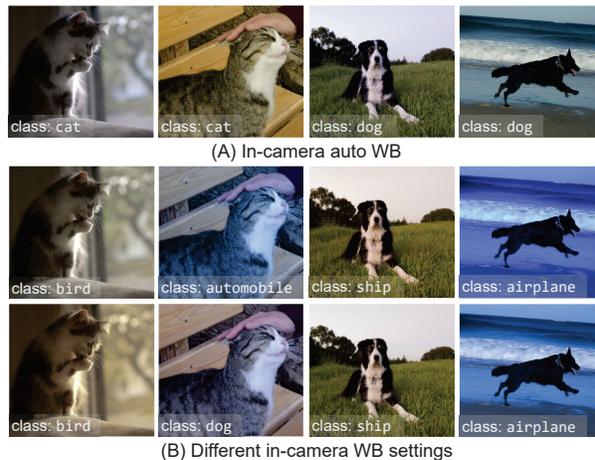


Figure 8. (A) Correctly classified images rendered with in-camera auto WB. (B) Misclassified images rendered with *in-camera* different WB. Note that all images in (B) are correctly classified by the same model (AlexNet [36]) trained on WB augmented data.

plemental materials.

## 6. Conclusion

This work has examined the impact on computational color constancy errors on DNNs for image classification and semantic segmentation. A new method to perform augmentation that accurately mimics WB errors was introduced. We show that both pre-processing WB correction and training DNNs with our augmented WB images improve the results for DNNs targeting CIFAR-10, CIFAR-100, and ADE20K datasets. We believe our WB augmentation method will be useful for other tasks targeted by DNN where image augmentation is sought.

**Acknowledgments** This study was funded in part by the Canada First Research Excellence Fund for the Vision: Science to Applications (VISTA) programme and an NSERC Discovery Grant. Dr. Brown contributed to this article in his personal capacity as a professor at York University. The views expressed are his own and do not necessarily represent the views of Samsung Research.

## References

- [1] Mahmoud Afifi and Michael S Brown. Sensor-independent illumination estimation for DNN models. In *BMVC*, 2019. 2
- [2] Mahmoud Afifi, Brian Price, Scott Cohen, and Michael S Brown. When color constancy goes wrong: Correcting improperly white-balanced images. In *CVPR*, 2019. 2, 4, 5, 6, 7, 8
- [3] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access*, 6:14410–14430, 2018. 2
- [4] Alexander Andreopoulos and John K Tsotsos. On sensor bias in experimental methods for comparing interest-point, saliency, and recognition algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(1):110–126, 2012. 3
- [5] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017. 6, 7, 8
- [6] Jonathan T Barron. Convolutional color constancy. In *ICCV*, 2015. 2
- [7] Jonathan T Barron and Yun-Ta Tsai. Fast fourier color constancy. In *CVPR*, 2017. 2
- [8] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. *arXiv preprint arXiv:1811.11127*, 2018. 2
- [9] Gershon Buchsbaum. A spatial processor model for object colour perception. *Journal of the Franklin Institute*, 310(1):1–26, 1980. 2
- [10] Gershon Buchsbaum. A spatial processor model for object colour perception. *Journal of the Franklin Institute*, 310(1):1–26, 1980. 5, 6
- [11] Vladimir Bychkovsky, Sylvain Paris, Eric Chan, and Frédo Durand. Learning photographic global tonal adjustment with a database of input/output image pairs. In *CVPR*, 2011. 7
- [12] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (SP)*, 2017. 2
- [13] Alexandra Carlson, Katherine A Skinner, and Matthew Johnson-Roberson. Modeling camera effects to improve deep vision for real and synthetic data. In *ECCV*, 2018. 3
- [14] A. Chakrabarti, Ying Xiong, Baochen Sun, T. Darrell, D. Scharstein, T. Zickler, and K. Saenko. Modeling radiometric uncertainty for vision with tone-mapped color images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(11):2185–2198, 2014. 2
- [15] Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*, 2014. 2
- [16] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 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, 2018. 3, 4
- [17] Dongliang Cheng, Brian Price, Scott Cohen, and Michael S Brown. Effective learning-based illuminant estimation using simple features. In *CVPR*, 2015. 2
- [18] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. In *ICML*, 2017. 2
- [19] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 2
- [20] Duc-Tien Dang-Nguyen, Cecilia Pasquini, Valentina Conotter, and Giulia Boato. Raise: A raw images dataset for digital image forensics. In *ACM Multimedia Systems Conference*, 2015. 7
- [21] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3, 4
- [22] Steven Diamond, Vincent Sitzmann, Stephen Boyd, Gordon Wetzstein, and Felix Heide. Dirty pixels: Optimizing image classification architectures for raw sensor data. *arXiv preprint arXiv:1701.06487*, 2017. 3
- [23] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2
- [24] Graham D Finlayson, Michal Mackiewicz, and Anya Hurlbert. Color correction using root-polynomial regression. *IEEE Transactions on Image Processing*, 24(5):1460–1470, 2015. 4
- [25] Graham D Finlayson and Elisabetta Trezzi. Shades of gray and colour constancy. In *Color and Imaging Conference*, 2004. 2
- [26] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 2
- [27] Samuel W Hasinoff, Dillon Sharlet, Ryan Geiss, Andrew Adams, Jonathan T Barron, Florian Kainz, Jiawen Chen, and Marc Levoy. Burst photography for high dynamic range and low-light imaging on mobile cameras. *ACM Transactions on Graphics*, 35(6):192, 2016. 2, 7
- [28] Søren Hauberg, Oren Freifeld, Anders Boesen Lindbo Larsen, John Fisher, and Lars Hansen. Dreaming more data: Class-dependent distributions over diffeomorphisms for learned data augmentation. In *Artificial Intelligence and Statistics*, 2016. 2
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 3, 4
- [30] Yuanming Hu, Baoyuan Wang, and Stephen Lin. Fc4: fully convolutional color constancy with confidence-weighted pooling. In *CVPR*, 2017. 2
- [31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 3
- [32] Nima Khademi Kalantari and Ravi Ramamoorthi. Deep high dynamic range imaging of dynamic scenes. *ACM Transactions on Graphics*, 36(4):144–1, 2017. 2

- [33] Hakki Can Karaimer and Michael S Brown. A software platform for manipulating the camera imaging pipeline. In *ECCV*, 2016. 2
- [34] Seon Joo Kim, Hai Ting Lin, Zheng Lu, Sabine Süssstrunk, Stephen Lin, and Michael S Brown. A new in-camera imaging model for color computer vision and its application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(12):2289–2302, 2012. 2
- [35] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009. 6, 7
- [36] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 2, 3, 4, 6, 7, 8
- [37] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *Technical report, Google, Inc.*, 2016. 2
- [38] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *ICCV*, 2017. 2
- [39] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*, 2017. 1, 3, 4
- [40] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *ICLR*, 2017. 2
- [41] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *CVPR*, 2016. 2
- [42] Yair Movshovitz-Attias, Takeo Kanade, and Yaser Sheikh. How useful is photo-realistic rendering for visual learning? In *ECCV*, 2016. 2
- [43] Hyeonseob Nam and Hyo-Eun Kim. Batch-instance normalization for adaptively style-invariant neural networks. In *NIPS*, 2018. 3
- [44] Seonghyeon Nam and Seon Joo Kim. Modelling the scene dependent imaging in cameras with a deep neural network. In *ICCV*, 2017. 2
- [45] Rang MH Nguyen and Michael S Brown. Raw image reconstruction using a self-contained sRGB–JPEG image with small memory overhead. *International Journal of Computer Vision*, 126(6):637–650, 2018. 2
- [46] Luis Perez and Jason Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017. 2, 6, 7, 8
- [47] Rajeev Ramanath, Wesley E Snyder, Youngjun Yoo, and Mark S Drew. Color image processing pipeline. *IEEE Signal Processing Magazine*, 22(1):34–43, 2005. 2
- [48] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 2
- [49] Alexander Jung Revision. Imgaug library. Online; accessed 30 January 2019. 2
- [50] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *ECCV*. 3
- [51] Wu Shi, Chen Change Loy, and Xiaoou Tang. Deep specialized network for illuminant estimation. In *ECCV*, 2016. 2
- [52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3, 4
- [53] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 3, 4
- [54] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014. 2
- [55] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*, 2017. 3
- [56] VSR Veeravasarapu, Constantin Rothkopf, and Ramesh Visvanathan. Adversarially tuned scene generation. In *CVPR*, 2017. 2
- [57] Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Song. Spatially transformed adversarial examples. *arXiv preprint arXiv:1801.02612*, 2018. 2
- [58] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples for semantic segmentation and object detection. In *ICCV*, 2017. 2
- [59] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2015. 3, 4
- [60] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. *arXiv preprint arXiv:1708.04896*, 2017. 2
- [61] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ADE20K dataset. In *CVPR*, 2017. 3, 4, 6, 7, 8