# Deep Photo Scan: Semi-Supervised Learning for dealing with the real-world degradation in Smartphone Photo Scanning

Man M. Ho
Hosei University
Tokyo, Japan
man.hominh.6m@stu.hosei.ac.jp

Jinjia Zhou
Hosei University
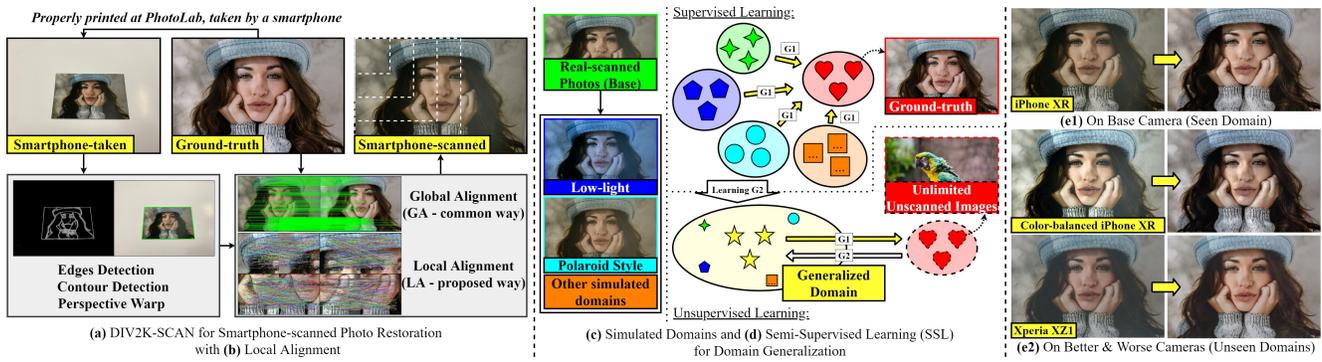Tokyo, Japan
jinjia.zhou.35@hosei.ac.jp

Figure 1. We present DIV2K-SCAN dataset for smartphone-scanned photo restoration (a) with Local Alignment (b), simulate varied domains to gain generalization in scanned image properties using low-level image transformation (c), and design a Semi-Supervised Learning system to train our network on also unscanned images, diversifying training image content (d). As a result, this work obtains state-of-the-art performance on smartphone-scanned photos in seen and unseen domains (e1-e2).

## Abstract

*Physical photographs now can be conveniently scanned by smartphones and stored forever as a digital version, yet the scanned photos are not restored well. One solution is to train a supervised deep neural network on many digital photos and the corresponding scanned photos. However, it requires a high labor cost, leading to limited training data. Previous works create training pairs by simulating degradation using image processing techniques. Their synthetic images are formed with perfectly scanned photos in latent space. Even so, the real-world degradation in smartphone photo scanning remains unsolved since it is more complicated due to lens defocus, lighting conditions, losing details via printing. Besides, locally structural misalignment still occurs in data due to distorted shapes captured in a 3-D world, reducing restoration performance and the reliability of the quantitative evaluation. To solve these problems, we propose a semi-supervised Deep Photo Scan (DPScan). First, we present a way of producing real-world degradation and provide the DIV2K-SCAN dataset for smartphone-scanned photo restoration. Also, Local Alignment is proposed to reduce the minor misalignment remaining in data. Second, we simulate many different variants of the real-world degradation using low-level image transformation to gain a generalization in smartphone-scanned image properties, then train a degradation network to generalize all styles of degradation and provide pseudo-scanned photos for unscanned images as if they were scanned by a smartphone. Finally, we propose a Semi-Supervised Learning that allows our restoration network to be trained on both scanned and unscanned images, diversifying training image content. As a result, the proposed DPScan quantitatively and qualitatively outperforms its baseline architecture, state-of-the-art academic research, and industrial products in smartphone photo scanning.*

## 1. Introduction

Every moment passing by is precious; especially, when it marks important life events such as graduation, wedding. The moments are usually captured in photographs. How-

| iPhone XR (base) | Xperia XZ1 | Taken at Dusk, Low-Light | A Polaroid Style |

Figure 2. We manually simulate different domains of degradation affected by other shooting environments (*taken at dusk with lack of light*) and devices (*Xperia XZ1*, *Polaroid Camera*) using low-level image transformation based on the real-world degradation. As a result, the simulated photos are qualitatively closed to the real-scanned photos, proving the feasibility of our approach in providing pseudo-scanned images for unscanned photos as if they were also taken in/by other shooting environments and devices.

ever, we do not always have a digital version since digital cameras were not common in the past, or we accidentally lost it. Thanks to technological development, photographs can be efficiently stored by scanning applications on smartphones as high-resolution digital images. It also provides an efficient way to share the captured moments in physical photographs to everyone through the internet. To restore the scanned photos, the recent Old Photo Restoration (OPR) [31, 32] based on Deep Neural Networks (DNNs) tries to mimic and learn the scanning degradation using low-level image processing techniques such as Gaussian Noise, Gaussian Blur. Their artificially degraded photos are then formed with the perfectly scanned old photos in latent space. However, the smartphone-scanned photos are still not restored well since they contain more complicated degradation caused by camera quality, scanning environments, losing details via printing, various post-processing techniques, etc. In this work, we adopt DIV2K [30] to present the DIV2K-SCAN dataset, which provides real-world degradation in smartphone photo scanning. Besides the common way of globally aligning a scanned photo to their ground-truth, we propose Local Alignment (LA) to reduce a minor misalignment remaining in data. Inspired by [10], based on the captured real-world degradation, we simulate many different variants using low-level image transformation to gain a generalization in smartphone-scanned image properties for our restoration network. Furthermore, we leverage the concept of Generative Adversarial Networks (GANs) [8, 14, 44] to first generalize all domains of degradation, then provide pseudo inputs for an unlimited amount of unscanned images in training. Being joint with supervised training, we design a cycle process as high-quality images → scanned photos/pseudo inputs → reconstructed images. The proposed semi-supervised scheme balances two supervised and unsupervised errors while optimizing to limit the effect of imperfect pseudo inputs but still enhance restoration. Our approach is briefly described in Figure 1. Besides, our code and data are available at `https://minhmanho.github.io/dpscan/`.

**Creating real-world degradation.** The performance of DNNs mostly depends on how the training data is created. Therefore, it is crucial to create a specific degradation that is close to real-world problems. For example, [30, 34, 6, 16] use traditional interpolation techniques to achieve the distortion representing the problem of super-resolution. However, they are limited in digital zooming. [42] thus presented a way to obtain ground-truth for zoomed regions by optically zooming and provide a dataset for real-world computational zoom. Their work thus outperforms the state-of-the-art super-resolution work [34] in the field. [24, 1] have the same motivation in solving the natural noise of photographs caused by low-ISO, [13] enhancing smartphone-taken photos by creating training {smartphone, DSLR camera} pairs. This work is the first time to have a dataset (DIV2K-SCAN) that provides real-world degradation for smartphone-scanned photo restoration. Besides, inspired by [10], we apply low-level image transformation to simulate varied variants of the captured degradation as if the photos were also captured in/by other shooting environments and smartphones for training. To demonstrate the feasibility of this scheme, we have provided a proof of concept shown in Figure 2. Our trained restoration network is thus generalized in smartphone-scanned image properties.

**Semi-Supervised Learning.** Since the human-annotated data costs a considerable resource, worldwide researchers have proposed many semi-supervised learning techniques to solve the lack of labeled data (ground-truth) while training a deep neural network. For example, the image classification works [3, 37, 5] utilized a well-trained model (teacher) to generate pseudo labels for unlabeled images so that the smaller network (student) can be trained on them. The semantic segmentation work [23] added an unsupervised loss on synthetic maps predicted by auxiliary decoders for unlabeled samples. They all provided good schemes when the ground-truth data is limited. However, in this work, we lack the input data for training. Thanks to the seminal work of Generative Adversarial Networks (GANs) [8], the synthetic images now are high-fidelity and closer to

**a) Supervised Learning:**
800 Ground-truth images (DIV2K-SCAN) have its scanned version

**b) Unsupervised Learning:**
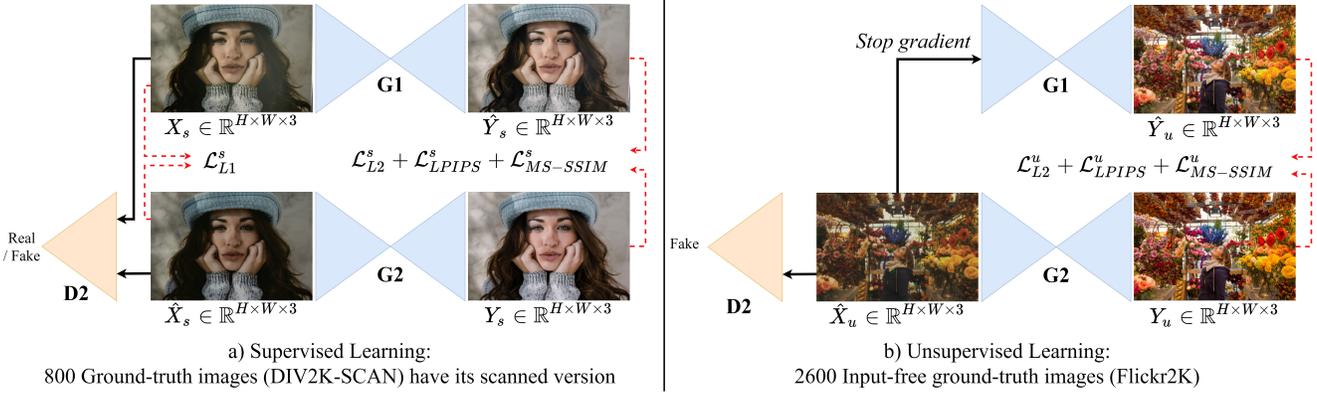2600 Input-free ground-truth images (Flickr2K)

Figure 3. We present a Semi-Supervised Learning system that allows our model to be trained on scanned (supervised) (a) and unscanned (unsupervised) (b) photos under strongly supervised loss functions such as $L2$ ($\mathcal{L}^*_{L2}$), LPIPS [40] ($\mathcal{L}^*_{LPIPS}$), and MS-SSIM [36] ($\mathcal{L}^*_{MS-SSIM}$), where $*$ denotes $s$ or $u$ representing supervised or unsupervised scheme respectively. Meanwhile, the distribution of the real-scanned photos is captured with adversarial losses and $L1$ ($\mathcal{L}^s_{L1}$). Errors between (a) and (b) are balanced during optimization.

the target distribution [14]. Moreover, the synthetic images can be translated back to the input domain, creating a cycle consistency [44, 35] for the task. Inspired by the aforementioned works, we adopt the concept of GANs to generalize our real-world degradation and its simulated versions to map an unlimited amount of unscanned images to a generalized domain, providing them the pseudo-scanned inputs for unsupervised training. Thus, the training image content is diversified. Being joint with supervised training, it will create a cycle fashion as high-quality images → scanned photos/pseudo inputs → reconstructed images, as shown in Figures 1 and 3. The proposed semi-supervised scheme balances the errors between supervised and unsupervised while optimizing to limit the effect of imperfect pseudo inputs but still enhance restoration.

**Network architecture**. The autoencoder architecture U-Net [27] and its variants have gained high performance and become well-known in image-to-image translation tasks [2, 41, 14, 11]. Recently, many techniques for customizing a deep neural network have also grown rapidly towards enhancing efficiency and effectiveness. It is meaningful for this smartphone photo scanning to operate on a limited resource. Inspired by the attention mechanism, Wang et al. [33] proposed an Efficient Channel Attention (ECA) module that reduces a huge computational cost with favorable performance compared to its backbones in image classification, object detection, and semantic segmentation. Therefore, we leverage ECA to design a Residual ECA (RECA) Block and RECA U-Block for our restoration network. As a consequence, the customized architecture shows learning capability compared with its baseline architecture (named as Simple DPScan) built by U-Net [27], residual modules [15] between encoder and decoder, blur pooling [39], and EvoNorm-S0 [20].

Our contributions are as follows:

- We present the DIV2K-SCAN dataset, which provides real-world degradation in smartphone photo scanning. Also, the dataset allows a deep neural network to be trained with strongly supervised loss functions.

- Although the smartphone-scanned image pairs are aligned globally, a minor misalignment still occurs, lowering restoration performance and making a quantitative comparison using similarity metrics less reliable. To address this problem, we propose a Local Alignment (LA) to perfectly align a smartphone-scanned photo to its ground-truth. LA-ed data also shows that the larger the image size, the more serious misalignment.

- To address the concern of our performance on photos captured in/by other shooting environments and devices (generalization), inspired by [10], we apply color style transfer to simulate varied types of degradation based on the real-world degradation. Thus, our work gains a generalization in smartphone-scanned image properties.

- We propose the semi-supervised Deep Photo Scan (DPScan) that has two advantages: a) Semi-Supervised Learning diversifying training image content by allowing our restoration network to be trained on both scanned and unscanned images and b) the customized Residual Efficient Channel Attention (RECA) Block and RECA U-Block. As a result, our semi-supervised DPScan outperforms its baseline, the previous research works, and industrial products comprehensively in 1-domain and generalization tests.

## 2. The Proposed Deep Photo Scan (DPScan)

Our work consists of two main components: (1) data preparation such as reproducing real-world degradation, image annotation/alignment, and simulating many smartphone scanning styles using color style transfer, and (2) the proposed semi-supervised DPScan for smartphone-scanned photo restoration.

Regarding (1), we leverage both traditional Canny [4] and DNN-based [25] edge detection techniques to identify the contour of interests. Afterward, the annotated images are warped and cropped to have a top-down view as though a professional scanner scanned them. Furthermore, we apply a precise alignment based on SIFT [21] and RANSAC [7, 29] to suppress the structural mismatch between inputs and ground-truth images with Local Alignment (LA). Besides, inspired by [10] and the proof shown in Figure 2, we apply color style transfer to sample many different domains based on the real-world degradation as if the photos were also scanned in/by other shooting environments and devices, as shown in Figure 1. Regarding (2), we design a semi-supervised framework for our DPScan that includes two generators $G1$ and $G2$, and a discriminator $D2$. In that, $G1$ is to restore scanned photos and trained under supervised loss functions. Meanwhile, the GAN-based $G2$ provides the pseudo inputs for unscanned images in a generalized domain and is trained with the discriminator $D2$, which can distinguish whether a scanned photo is real or fake. Initially, all models are pre-trained on 1-domain DIV2K-SCAN (iPhone XR) with a supervised learning scheme, that $G1$ is trained independently with $G2$ and $D2$. Afterward, $G1$, $G2$, and $D2$ are jointly trained on scanned photos from DIV2K-SCAN and unscanned images from Flick2K [30], representing a Semi-Supervised Learning for dealing with the lack of inputs. In case of being fine-tuned on multiple-domain DIV2K-SCAN, $G2$ will generalize all domains and provide pseudo-scanned photos in the generalized domain, as shown in Figures 1 and 3. **Please check our supplemental document for a visualization of pseudo-scanned photos**.

### 2.1. Supervised Learning for pre-training G1, G2 and D2

In pre-training on DIV2K-SCAN with supervised learning, after perspective warping, our $G1 : X \rightarrow Y$ restores the scanned inputs $X_s \in \mathbb{R}^{H \times W \times 3}$ to have $\hat{Y}_s \in \mathbb{R}^{H \times W \times 3}$, as follows:

$$\hat{Y}_s = G1(X_s) \tag{1}$$

The errors between $\hat{Y}_s \in \mathbb{R}^{H \times W \times 3}$ and its ground-truth images $Y$ are optimized under several supervised losses such as $L2$, Multiscale Structural Similarity (MS-SSIM) [36], and the perceptual metric LPIPS [40] as follows:

$$\mathcal{L}_{G1}^s = \alpha * \mathcal{L}_{L2}^s + \beta * \mathcal{L}_{LPIPS}^s + \gamma * \mathcal{L}_{MS\text{-}SSIM}^s \tag{2}$$

where:
$\mathcal{L}_{L2}^s = ||Y_s - \hat{Y}_s||_2^2$,
$\mathcal{L}_{MS\text{-}SSIM}^s = MS\text{-}SSIM(Y_s, \hat{Y}_s)$ described in [36],
$\mathcal{L}_{LPIPS}^s = LPIPS(Y_s, \hat{Y}_s)$ described in [40].
To learn the specific degradation of scanned photos from DIV2K-SCAN, we independently train a simply-designed network $G2 : Y \rightarrow X$ to degrade the ground-truth images $Y_s$ to synthesize its scanned version $\hat{X}_s \in \mathbb{R}^{H \times W \times 3}$, as follows:

$$\hat{X}_s = G2(Y_s) \tag{3}$$

Instead of conditioning the discriminator [22, 38], we train the generator $G2$ and discriminator $D2$ under the $L1$ [14] and hinge adversarial losses [18] as:

$$\mathcal{L}_{D2}^s = -\mathbb{E}_{X_s}[min(0, -1 + D(X_s)] \\ - \mathbb{E}_{Y_s}[min(0, -1 - D(G2(Y_s))] \tag{4}$$

$$\mathcal{L}_{G2}^s = -\mathbb{E}_{Y_s}[D(G2(Y_s))] \tag{5}$$

$$\mathcal{L}_{L1}^s = \mathbb{E}_{X_s, Y_s}[||X_s - G2(Y_s)||_1] \tag{6}$$

The total loss for G2 is defined as:

$$\mathcal{L}_{G2\_final}^s = \alpha * \mathcal{L}_{L1}^s + \delta * \mathcal{L}_{G2}^s \tag{7}$$

We empirically set $\alpha = 1$, $\beta = 0.2$, and $\gamma = 1$, $\delta = 0.05$ for this supervised learning scheme.

### 2.2. Semi-Supervised Learning for fine-tuning G1, G2, and D2 together

After pre-training models on DIV2K-SCAN, we then train $G1$, $G2$, and $D2$ together on both DIV2K-SCAN providing the ground-truth images $Y_s$ with their scanned photos $X_s$ and Flick2K [30] providing the input-free ground-truth images $Y_u \in \mathbb{R}^{H \times W \times 3}$.

Firstly, $X_s$ and $Y_s$ are processed as described in Section 2.1. Secondly, $G2$ provides pseudo inputs $\hat{X}_u$ for $Y_u$ as:

$$\hat{X}_u = G2(Y_u) \tag{8}$$

under adversarial losses updated from Equations 4 and 5 as:

$$\mathcal{L}_{D2} = -\mathbb{E}_{X_s}[min(0, -1 + D(X_s)] \\ - 0.5 * (\mathbb{E}_{Y_s}[min(0, -1 - D(G2(Y_s))] \\ + \mathbb{E}_{Y_u}[min(0, -1 - D(G2(Y_u))]) \tag{9}$$

$X \in \mathbb{R}^{H \times W \times 3}$   32  32  64  128  256  512  512  512  512  256  128  64  32  32  3   $\hat{Y} \in \mathbb{R}^{H \times W \times 3}$

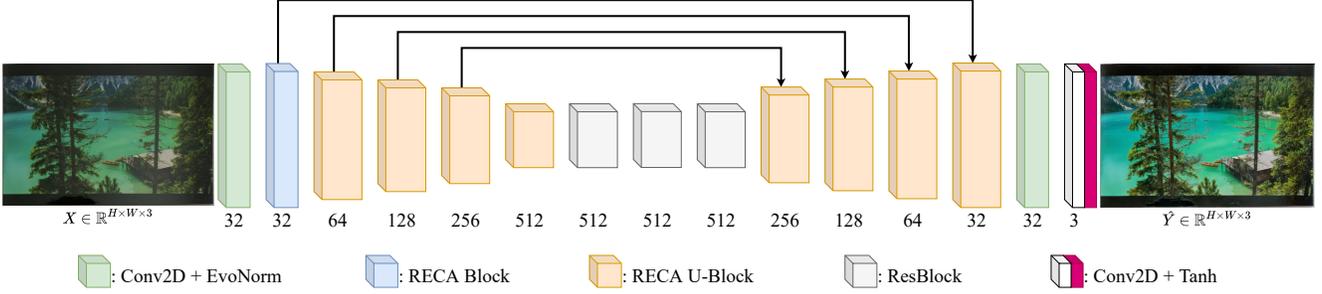: Conv2D + EvoNorm    : RECA Block    : RECA U-Block    : ResBlock    : Conv2D + Tanh

Figure 4. Architecture of our restoration network $G1$ restoring the scanned photo $X$ to have its high-quality $\hat{Y}$.
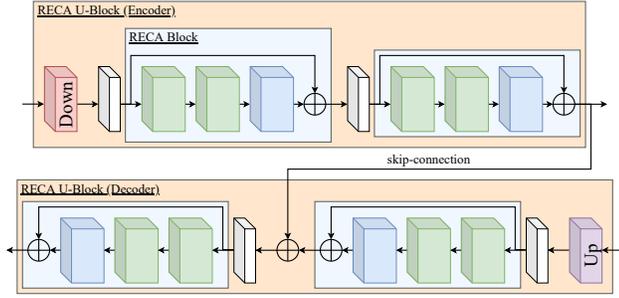


Figure 5. Illustrations of Residual Efficient Channel Attention (RECA) Block and RECA U-Block. Regarding "*DOWN*"-/"*UP*"-sampling, we use the anti-aliasing max pooling and bi-linear interpolation from [39]. $\oplus$ represents a summation.

$$\mathcal{L}_{G2} = -0.5 * \left( \mathbb{E}_{Y_s}[D(G2(Y_s))] \\ + \mathbb{E}_{Y_u}[D(G2(Y_u))] \right) \quad (10)$$

Combined with the supervised loss in Equation 6, we have the updated total loss for $G2$ as:

$$\mathcal{L}_{G2\_final} = \alpha * \mathcal{L}_{L2}^s + \delta * \mathcal{L}_{G2} \quad (11)$$

Afterward, $G1$ leverages the pseudo input $\hat{X}_u$ to synthesize the reconstructed $\hat{Y}_u$ creating a cycle fashion as:

$$\hat{Y}_u = G1(sg(\hat{X}_u)) = G1(sg(G2(Y_u))) \quad (12)$$

where $sg$ denotes the $stop\_gradient$ function added to avoid falsifying the distribution of real-scanned photos. Similar to Equation 2, the defined loss function for $G1$ optimizing errors between $\hat{Y}_u$ and $Y_u$ is as:

$$\mathcal{L}_{G1}^u = \alpha * \mathcal{L}_{L2}^u + \beta * \mathcal{L}_{LPIPS}^u + \gamma * \mathcal{L}_{MS-SSIM}^u \quad (13)$$

Finally, the semi-supervised loss function for $G1$ is as:

$$\mathcal{L}_{G1} = \eta * \mathcal{L}_{G1}^s + (1 - \eta) * \mathcal{L}_{G1}^u \quad (14)$$

where $\eta$ is a balance weight between supervised and unsupervised errors, and is set to $0.5$. Other hyper-parameters empirically re-set $\alpha = 1$, $\beta = 0.1$, and $\gamma = 0.25$, $\delta = 0.05$ for this Semi-Supervised Learning scheme.

## 2.3. Network architecture

The proposed semi-supervised DPScan consists of three main deep neural networks: $G1 : X \rightarrow Y$ for scanned photo restoration, $G2 : Y \rightarrow X$ for degrading high-quality images, and discriminator $D2$ trained together with $G2$ to distinguish whether the scanned images are real or fake.

**Generator** $G1$ is designed based on the network architecture of U-Net [27] with skip connections [11], residual modules (ResBlock) between the encoder and decoder [9, 15], EvoNorm-S0 [20], which have achieved a high performance in image-to-image translation tasks. Besides, in each block of the encoder and decoder, we leverage the anti-aliasing max pooling and bi-linear interpolation [39] for our down-sampling and up-sampling, respectively. The network architecture adopting the mentioned techniques is our baseline architecture, named Simple DPScan. Afterward, we adopt Efficient Channel Attention (ECA) module [33], which has shown the efficiency but effectiveness in image classification, to design Residual ECA (RECA) and RECA U-Block and customize $G1$, as described in Figures 4 and 5. Please find technical details in our supplemental document.

**Generator** $G2$ **and Discriminator** $D2$. We utilize the baseline architecture mentioned above for $G2$ to generate pseudo-scanned photos. Meanwhile, $D2$ is entirely based on the discriminator of SA-GAN [38] with Spectral Normalization [22]. Please find more details in our supplemental document.

## 2.4. Data preparation and DIV2K-SCAN

**Generating the real-world degradation**. We first rotate the high-quality images from DIV2K [30] so that all images are in landscape format, then center crop all images to an aspect ratio of $15 : 10$. Afterward, we ask a professional

*Sliding Window for Patch Extraction*

Globally-Aligned Pair

Center Crop $R_1$ %
Resize to $M \times N$

Simplest Color Balance

Extract Patches Patch Size of $W_1$

Color-balanced    Ground-truth

Find Homography SIFT + RANSAC

Perspective Warp

Center Crop $R_2$ %
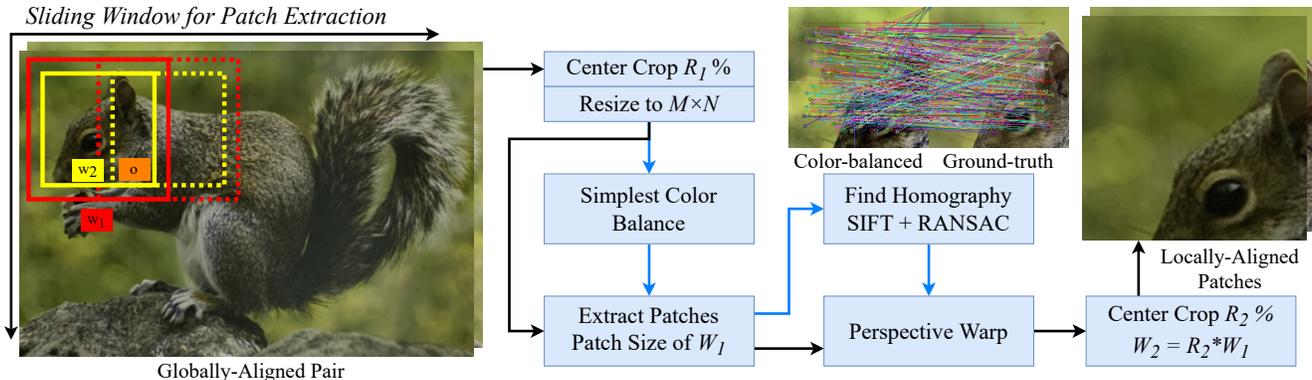$W_2 = R_2 * W_1$

Locally-Aligned Patches

Figure 6. Generating Locally-Aligned data. After globally warping, we step-by-step apply a center crop to $R_1\%$ of the current size to remove the black borders, resize to $M \times N$ using bicubic interpolation, extract patches from color-balanced [19] photos to find homography matrices (blue) and from original photos for warping using a sliding window with a size of $W_1$ and a stride of $S\%$ of $W_1$, warp the scanned patches, center crop to $R_2\%$ of the size again, and finally obtain the locally-aligned patches with a size of $W_2 = R_2 * W_1$. $O = 1 - S/R_2$ denotes the percentage of how much two consecutive final patches overlap. Extracting and warping patches are powered by Kornia [26].

photography lab, where the staff is well-trained to print photographs with accurate colors and high quality, to print the processed images out with a size of $7.5cm \times 5cm$. All physical photos are then digitally taken in a room with sufficient light intensity on the white background using a smartphone. Therefore, our generated data contains a complex degradation of smartphone-scanned photos such as natural noises, haze, smartphone-level image quality, structural distortion, lost details via printing, etc.

**Contour detection and image alignment**. After collecting the digital images of printed photos, we apply the efficient Canny [4] edge detection to detect the contours of the actual scanned photos in the white background. However, the method is sensitive to the color between the boundary and usually makes mistakes in detecting a complete contour. Thus, we utilize learning-based DexiNed [25] to support Canny's method. The remaining failed contour detection will be manually corrected by humans. **Please check our supplemental document for discussion and visualization**.

Even though the top-down view is obtained, the structural mismatch between the warped image and its ground-truth still occurs. It becomes more challenging to train a deep neural network. The recent work RANSAC-flow [28] presents an advanced technique to precisely-structurally align an image to another one having the same context; however, the degradation of smartphone-scanned images can be falsified by their warping. To reduce structural mismatch while keeping the degradation intact, we leverage the SIFT features [21] combined with RANSAC [7, 29] to find the homography, then align the warped images to their ground-truth. We eventually achieve 900 ground-truth images with a real-scanned version for supervised training, validation, and test. Unfortunately, the local misalignment

still occurs because the object shape of the photo taken in the 3-D world is usually distorted. Consequently, learning capability is harmed, and the quantitative evaluation using similarity metrics becomes less reliable. To address this issue, we present a way of generating Locally-Aligned (LA) data for training and test, as described in Figure 6.

**Training data**. We utilize 800 scanned image pairs from DIV2K-SCAN and $2,600$ unscanned images from Flickr2K [17] to extract $12,000$ pairs and $20,800$ unscanned patches, respectively, using the presented Local Alignment (described in Figure 6) with $M \times N$ of $1080 \times 720$, $R_1$ of $95\%$, $R_2$ of $95\%$, final patch size $W_2$ of $256 \times 256$, the stride $S$ of $65\%$ $W_2$ resulting in two consecutive final patches overlapping $O = 1 - S/R_2 \approx 31.57\%$. Based on the real-world degradation, we synthesize more $K * 12,000$ scanned photos as if they were also scanned in/by other environments and devices with a $K$ of 100 color styles collected by [10]. In training, the data is augmented by a random flip in horizontal and vertical ways and random rotation with the degrees of $0, 90, 180, 270$.

**Validation and test data**. From 100 globally-aligned images from a domain DIV2K-SCAN, we apply Local Alignment to extract $4000, 1500, 600, 100, 100$ patches with the final patch size $W_2$ of $176 \times 176$, $256 \times 256$, $384 \times 384$, $576 \times 576$, $1072 \times 720$, respectively, $M \times N$ of $1072 \times 720$, $R_1$ of $95\%$, $R_2$ of $80\%$, the stride $S$ of $50\%$ $W_2$ resulting in two consecutive final patches overlapping $O = 1 - S/R_2 \approx 37.5\%$. Regarding the final patch size $W_2$ of $1072 \times 720$, we only apply the first center crop and resize the photos. The pairs from each size is split into two sets, $40\%$ for validation (*valset*) and $60\%$ for test (*testset*). In each set, *valset* contains degradation styles of iPhone XR and two simulated domains unseen from training, and *testset* contains photos of iPhone XR, Color-Balanced [19]

6

Figure 7. Ablation study on RECA for customizing DPScan, training on Locally-Aligned (LA) data, and the proposed Semi-Supervised Learning (SSL). As a result, each component has an improvement in restoring the edges and reducing artifacts. **Check our supplemental video for a better comparison**.

| Method | Learning | Alignment | RECA | PSNR↑ | LPIPS↓ | MS-SSIM↑ |
|--------|----------|-----------|------|-------|--------|----------|
| Pix2Pix [14] | SL | GA | | 22.63 | 0.2138 | 0.8979 |
| CycleGAN [44] | SL | GA | | 20.24 | 0.2504 | 0.8836 |
| 1D-DPScan | SL | GA | | 23.78 | 0.1606 | 0.9275 |
| | SL | GA | ✓ | 24.10 | 0.1424 | 0.9316 |
| | SSL | GA | ✓ | 24.38 | 0.1423 | 0.9333 |
| | SL | LA | ✓ | 24.85 | 0.1351 | 0.9415 |
| | SSL | LA | ✓ | **25.26** | **0.1242** | **0.9446** |

Table 1. Comparison between previous works [14, 44] and our ablation models using (a) Supervised Learning (SL) or Semi-Supervised Learning (SSL), (b) Global Alignment (GA) or the proposed Local Alignment (LA), and (c) RECA. All models are trained and evaluated on 1-domain DIV2K-SCAN. ↑ / ↓: higher/lower is better.

iPhone XR, and Sony Xperia XZ1. In a comparison with industrial products, we set a $R_1$ of $85\%$ to avoid large black borders produced by Google Photo Scan.

## 3. Experiments

In this section, we conduct an ablation study on our presented techniques, including RECA for DPScan, training on Locally-Aligned (LA-ed) data, and the proposed Semi-Supervised Learning (SSL). Furthermore, we also train and compare with two typical works **Pix2Pix** [14] and **Cycle-GAN** [44] in the same condition on 1-domain DIV2K-SCAN (iPhone XR). DPScan trained on only iPhone XR is denoted as 1D-DPScan. Besides, to prove our generalization performance, we compare Generalized DPScan (G-DPScan) with the recent works that aim to solve many different degradation types of scanned photos such as **Industrial products Google Photo Scan (GPS) and Genius Scan (GS)** (we manually produce their results using iPhone XR), and **Academic research Old Photo Restoration** [32] on also unseen sets of color-balanced [19] iPhone XR and Xperia XZ1, which have better and worse performance than iPhone XR, respectively. Other experiments on deep modules for designing DPScan, image alignment,

pseudo-scanned photo synthesis, the number of simulated domains $K$ can be found in the supplemental document.

All quantitative comparisons are conducted on DIV2K-SCAN testsets described in Section 2.4 using similarity metrics such as Peak Signal-to-Noise Ratio (PSNR), LPIPS [40], and MS-SSIM [36]. Unfortunately, the similarity metrics are less reliable due to a minor misalignment remaining in data. Concretely, the smaller the image size, the smaller the structural mismatch between input and ground-truth after alignment, the more accurate the similar metrics, as shown in Figure 9. Therefore, we use the average score over three sizes of $176 \times 176$, $256 \times 256$, $384 \times 384$ for all quantitative comparisons on LA-ed data.

**Comparison between ablation models and previous works trained and evaluated on 1-domain DIV2K-SCAN (iPhone XR).** We adopt U-Net [27], residual modules [9, 15], anti-aliasing down-/up-samplers [39], EvoNorm [20] to design a baseline architecture, named Simple DPScan. While considering improving network architecture, we conduct an ablation study on customized deep learning techniques such as Flow Warping Block (FWB) [12], Residual Feature-based Attention (RFA), Residual Self-Attention (RSA) [38], Residual Channel Attention Block (RCAB) [43], and Residual Efficient Channel Attention (RECA) [33]. The Efficient Channel Attention (ECA) [33] has shown the efficiency yet effectiveness in reducing computational costs with high accuracy for the image classification task. Moreover, an experimental result shows that the customized RECA outperforms other ablation techniques with the best image quality (technical and experimental details are in the supplemental document). Therefore, we leverage the RECA module and its variant RECA U-Block, which is customized for u-style architecture [27], to design our network, as shown in Figures 4 and 5. We also compare our ablation models with the Pix2Pix [14] and CycleGAN [44] trained in the same condition to prove our restoration effectiveness. Besides, we propose Local Alignment and Semi-Supervised Learning (SSL) to solve (i) the remaining minor misalignment between the input and ground-truth in data and (ii) expensive costs leading to lack of real-scanned data. As a qualitative result, each presented component gradually improves the restoration performance as clearer edges with fewer artifacts, as shown in Figure 7 (a better qualitative comparison is in our supplemental video). In comparison with previous works, our work provides the highest image quality without the haze effect, as shown in Figure 10. Quantitatively, our designed baseline architecture for DPScan can outperform the previous works [14, 44] with better average PSNR, LPIPS [40], and MS-SSIM [36] of **23.78, 0.1606, 0.9275**. The average **PSNR** is further improved **+0.32dB** when the baseline DPScan is customized with RECA, **+0.75dB** more when the model is trained on LA-ed data, **+0.41dB** more when the
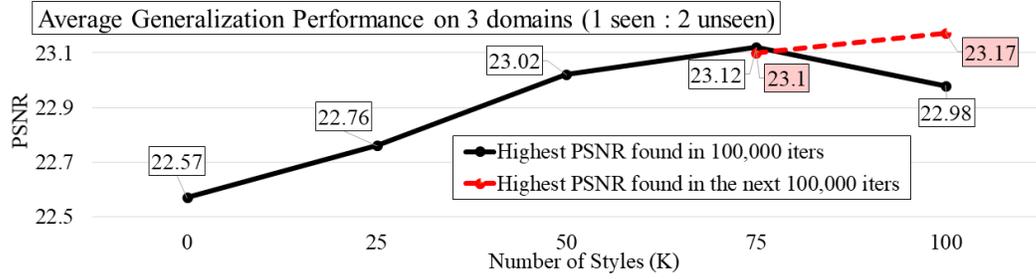
Figure 8. Ablation study on the number of simulated domains ($K$) for fine-tuning 1D-DPScan in two stages of $100,000$ iterations using average PSNR. The model with $K = 100$ obtains the highest generalization performance, even though it takes a longer training time. **Please check our supplemental document for more details**.
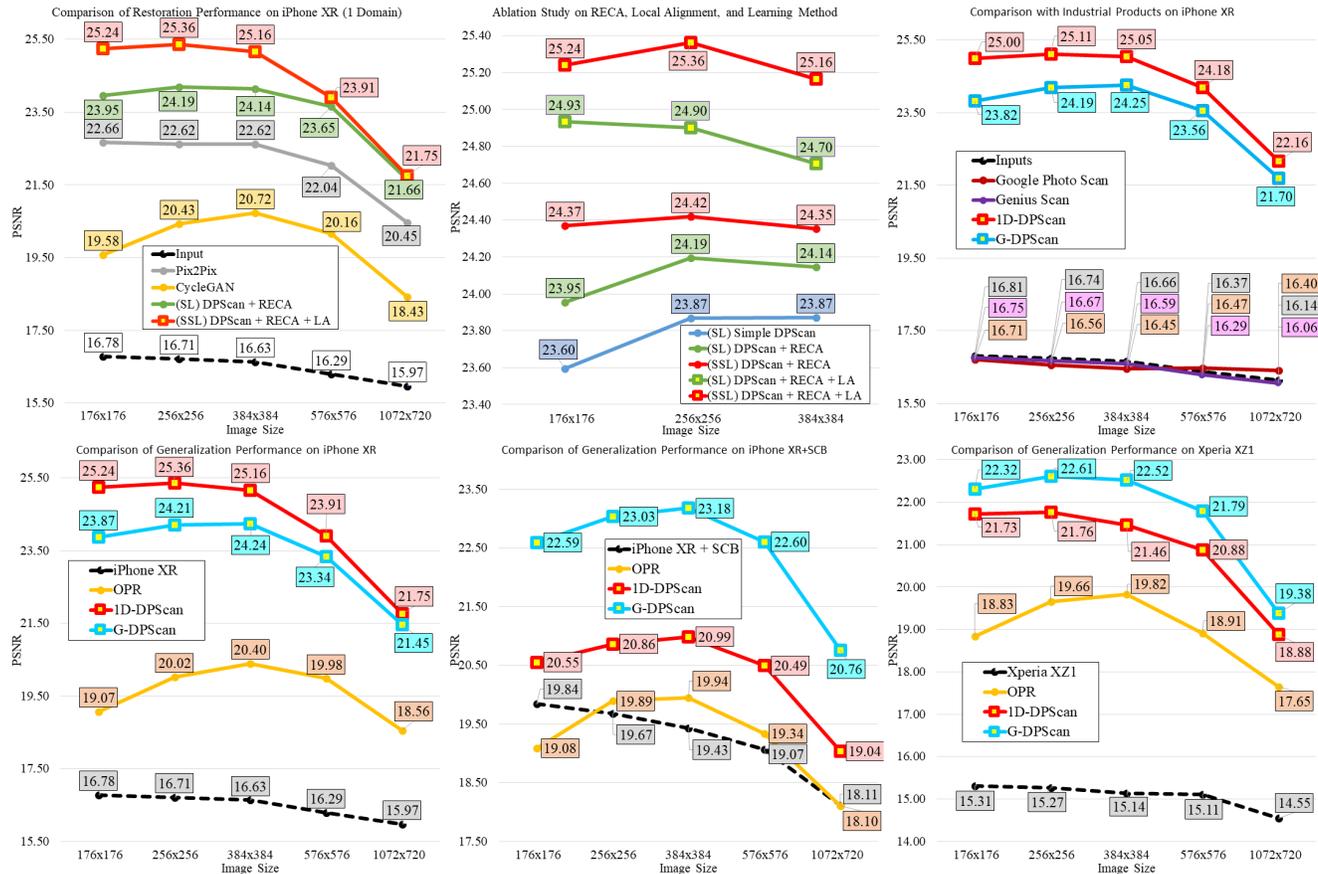


Figure 9. A full version of quantitative comparison using PSNR (*higher is better*) on multiple-domain DIV2K-SCAN (iPhone XR is a seen domain, while iPhone XR + SCB [19] and Xperia XZ1 are unseen domains) with an image size from $176 \times 176$ to $1072 \times 720$. Ablation models, 1D-DPScan, and the previous works Pix2Pix [14] and CycleGAN [44] are trained on and to solve 1-domain DIV2K-SCAN (iPhone XR); meanwhile, other methods such as Old Photo Restoration (OPR) [32], industrial products, and our G-DPScan are to solve multiple domains. This experiment shows that 1) the image quality is gradually reduced in ascending order of image size, proving that the larger the image size, the more serious misalignment, 2) each presented technique provides a significant improvement, and the final version of DPScan outperforms all ablation models (*middle-top*). 3) Our 1D-DPScan (trained on iPhone XR only) and G-DPScan (trained to solve multiple domains) outperform the research works [14, 44, 32] and industrial products Google Photo Scan and Genius Scan comprehensively. **Please check our supplemental document for a comparison using LPIPS and MS-SSIM**.

model is trained with the proposed SSL. Eventually, all presented techniques bring **+1.48dB** totally. Also, the average LPIPS and MS-SSIM are improved **-0.0364** and **+0.0171**, respectively, in total, as shown in Table 1 and Figure 9.

8

| Method | iPhone XR (seen) | | | iPhone XR + SCB [19] (unseen) | | | Xperia XZ1 (unseen) | | | Average | | | Method | iPhone XR (seen) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | LPIPS↓ | MS-SSIM↑ | PSNR↑ | LPIPS↓ | MS-SSIM↑ | PSNR↑ | LPIPS↓ | MS-SSIM↑ | PSNR↑ | LPIPS↓ | MS-SSIM↑ | | PSNR | LPIPS | MS-SSIM |
| Inputs | 16.71 | 0.3734 | 0.8436 | 19.65 | 0.3569 | 0.8853 | 15.24 | 0.4327 | 0.7783 | 17.20 | 0.3877 | 0.8357 | Inputs | 16.74 | 0.3944 | 0.8367 |
| OPR [32] | 19.83 | 0.3479 | 0.8792 | 19.64 | 0.3498 | 0.8786 | 19.44 | 0.4010 | 0.8551 | 19.64 | 0.3662 | 0.8710 | GPS | 16.58 | 0.3916 | 0.8431 |
| 1D-DPScan | **25.26** | **0.1242** | **0.9446** | 20.80 | 0.1883 | 0.9172 | 21.65 | 0.2357 | 0.8972 | 22.57 | 0.1827 | 0.9197 | GS | 16.67 | 0.3955 | 0.8322 |
| G-DPScan | 24.10 | 0.1413 | 0.9363 | **22.93** | **0.1610** | **0.9276** | **22.48** | **0.2134** | **0.9045** | **23.17** | **0.1719** | **0.9228** | G-DPScan | **24.09** | **0.1571** | **0.9315** |
| | a) w/ Recent Works. R1=95% (R1: The first center crop ratio set to remove black borders.) | | | | | | | | | | | | b) w/ Industrial Products. R1=85% | | | |

Table 2. A quantitative comparison of generalization performance (1 seen: 2 unseen domains) between the research work OPR [32], our DPScan trained on iPhone XR only (1D-DPScan), Generalized DPScan (G-DPScan), and industrial products Google Photo Scan (GPS) and Genius Scan (GS). ↑ / ↓: higher/lower is better.
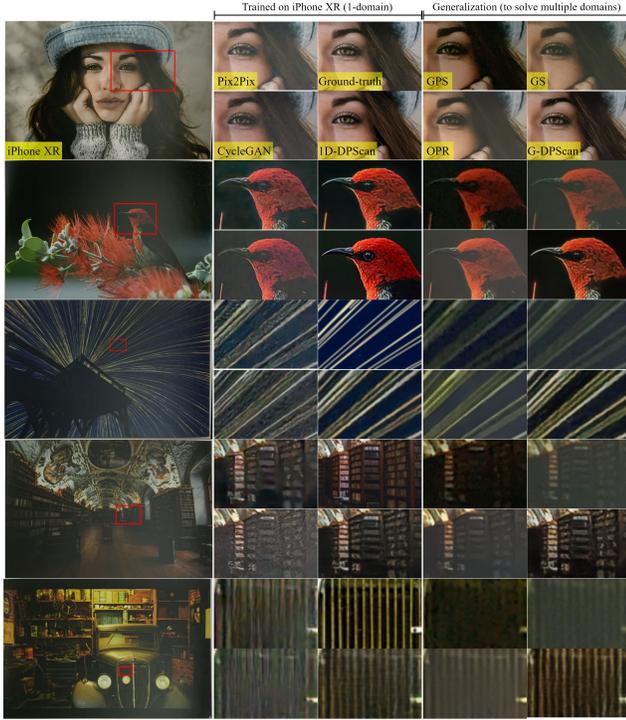


Figure 10. Qualitative comparison between two typical works Pix2Pix [14] and CycleGAN [44] trained on 1-domain DIV2K-SCAN (iPhone XR), industrial products Google Photo Scan (GPS) and Genius Scan (GS), the previous work Old Photo Restoration (OPR) [32], and our 1-domain (1D-DPScan) and generalized (G-DPScan) networks. This work provides the most detailed photos without haze and color fading. **Please check our supplemental document for a full version and more results**.
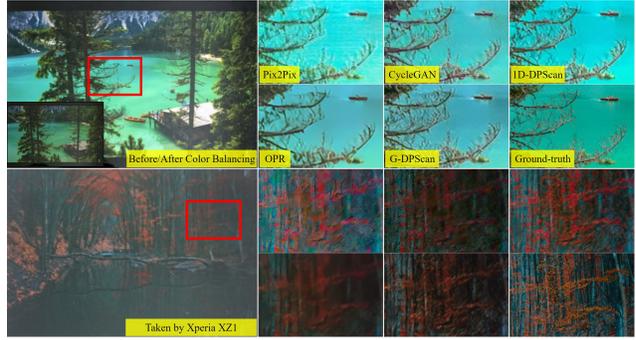


Figure 11. Qualitative comparison on unseen domains including the color-balanced [19] iPhone XR (*top sample*) and Xperia XZ1 (*bottom sample*). Our 1-domain (DPScan) and generalized (G-DPScan) models obtains the highest image quality. **Please check our supplemental document for a full version and more results**.

model is weakly trained because the ground-truth images of real-scanned photos are missing. Moreover, the real-world degradation in smartphone photo scanning is more complicated caused by real lens defocus, lighting conditions, many different smartphone post-processing techniques, etc.

To overcome the aforementioned issues, we present DIV2K-SCAN providing real-world degradation, Local Alignment effectively reducing the remaining structural misalignment in data, RECA-customized architecture, scanned photo degradation simulation (inspired by [10]) for domain generalization in smartphone-scanned image properties, and Semi-Supervised Learning diversifying training image content by allowing our network to be trained on scanned and unscanned images.

As a quantitative result, our 1-domain DPScan (1D-DPScan) can outperform OPR with better average **PSNR, LPIPS, MS-SSIM** of **22.57, 0.1827, 0.9197**. However, 1D-DPScan is trained on iPhone XR only, and its performance on iPhone XR is much higher than on unseen photos from Simplest-Color-Balanced [19] iPhone XR and Xperia XZ1, raising a concern of our generalization. Therefore, we simulate $K$ scanning domains based on our real-world degradation using color style transfer [10] as if our scanned photos were also taken in/by other environments and devices. An ablation study on $K \in \{25, 50, 75, 100\}$ shown in Figure 8 reveals that Generalized DPScan (G-DPScan) has the best

**Comparison with previous works and industrial products on multiple-domain DIV2K-SCAN**. The recent works Google Photo Scan (GPS), Genius Scan (GS), and Old Photo Restoration (OPR) [32] try to solve many different types of real-world degradation in smartphone photo scanning. Even though the industrial products GPS and GS have a user-friendly interface, their scanned photos are distorted with significant artifacts. Inspired by deep learning, OPR [32] trained their network on the scanning degradation simulated by low-level image processing techniques. Although their synthetic images are formed with the perfectly real-scanned old photos in latent space, their

generalization performance when $K = 100$, even though it takes a long training time (more details are in the supplemental document). Eventually, our G-DPScan gains better average **PSNR, LPIPS, MS-SSIM** of **23.17, 0.1780, 0.9223** compared with 1D-DPScan and the research work OPR [32], and **24.24, 0.1615, 0.9330** compared with industrial products GPS and GS, as shown in Table 2 and Figure 9.

Qualitatively, our DPScan provides the clearest edges without haze effect in both 1-domain and generalization tests. Concretely, our results show the highest details of the *girl*, *bird*, and *lines* on iPhone XR, as shown in Figure 10. Regarding the domains unseen from training and validation, such as Simplest-Color-Balanced [19] iPhone XR and Xperia XZ1, 1D-DPScan and G-DPScan provide the best image quality compared with the previous works. Surprisingly, 1D-DPScan generates more good-looking colors on its unseen domains than G-DPScan, although it is trained on iPhone XR only and quantitatively worse than G-DPScan, as shown in Figure 11. **Please check our supplemental document for more interesting experiments and results**. In conclusion, our semi-supervised DPScan outperforms the previous works and industrial products comprehensively.

## 4. Conclusion

We present a way to produce real-world degradation in smartphone photo scanning and present DIV2K-SCAN for smartphone-scanned photo restoration. Besides, Local Alignment is proposed to solve a minor misalignment remaining in data, which causes the reduction of restoration performance and reliability of the quantitative evaluation. Besides, we apply color style transfer to simulate many different variants of the real-world degradation as if the photos were also captured in/by other shooting environments and devices. Furthermore, we adopt the concept of GANs to degrade a high-quality image as if it were scanned by a smartphone and propose the semi-supervised Deep Photo Scan (DPScan) that has two advantages: 1) Semi-Supervised Learning allowing our network to be trained on both scanned and unscanned images and 2) the u-style architecture customized by Residual Efficient Channel Attention (RECA). Our work thus obtains a generalization in both training image content and smartphone-scanned image properties. As a result, our semi-supervised DPScan outperforms its baseline [27, 39, 20], the industrial Google Photo Scan and Genius Scan, the recent work Old Photo Restoration [31, 32], two retrained Pix2Pix [14] and CycleGAN [44] in 1-domain and generalization tests quantitatively and qualitatively. This work, therefore, becomes a promising baseline for smartphone-scanned photo restoration.

## References

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1692–1700, 2018.

[2] 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.

[3] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5049–5059, 2019.

[4] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, 1986.

[5] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020.

[6] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[7] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.

[8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[10] Man M. Ho and Jinjia Zhou. Deep preset: Blending and retouching photos with color style transfer. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2113–2121, January 2021.

[11] Minh-Man Ho, Jinjia Zhou, and Yibo Fan. Respecting low-level components of content with skip connections and semantic information in image style transfer. In *European Conference on Visual Media Production*, pages 1–9, 2019.

[12] Man M Ho, Jinjia Zhou, Gang He, Muchen Li, and Lei Li. Sr-cl-dmc: P-frame coding with super-resolution, color learning, and deep motion compensation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 124–125, 2020.

[13] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. Dslr-quality photos on mobile devices with deep convolutional networks. In *Proceedings*

*of the IEEE International Conference on Computer Vision*, pages 3277–3285, 2017.

[14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, 2017.

[15] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.

[16] Zhen Li, Jinglei Yang, Zheng Liu, Xiaomin Yang, Gwanggil Jeon, and Wei Wu. Feedback network for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[17] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.

[18] Jae Hyun Lim and Jong Chul Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017.

[19] Nicolas Limare, Jose-Luis Lisani, Jean-Michel Morel, Ana Belén Petro, and Catalina Sbert. Simplest color balance. *Image Processing On Line*, 1:297–315, 2011.

[20] Hanxiao Liu, Andrew Brock, Karen Simonyan, and Quoc V Le. Evolving normalization-activation layers. *arXiv preprint arXiv:2004.02967*, 2020.

[21] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.

[22] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.

[23] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020.

[24] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1586–1595, 2017.

[25] Xavier Soria Poma, Edgar Riba, and Angel Sappa. Dense extreme inception network: Towards a robust cnn model for edge detection. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1923–1932, 2020.

[26] Edgar Riba, Dmytro Mishkin, Daniel Ponsa, Ethan Rublee, and Gary Bradski. Kornia: an open source differentiable computer vision library for pytorch. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3674–3683, 2020.

[27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[28] Xi Shen, François Darmon, Alexei A Efros, and Mathieu Aubry. Ransac-flow: generic two-stage image alignment. In *16th European Conference on Computer Vision*, 2020.

[29] Richard Szeliski. Image alignment and stitching: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 2(1):1–104, 2006.

[30] Radu Timofte, Shuhang Gu, Jiqing Wu, Luc Van Gool, Lei Zhang, Ming-Hsuan Yang, Muhammad Haris, et al. Ntire 2018 challenge on single image super-resolution: Methods and results. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.

[31] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Bringing old photos back to life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2747–2757, 2020.

[32] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Old photo restoration via deep latent space translation. *arXiv preprint arXiv:2009.07047*, 2020.

[33] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11534–11542, 2020.

[34] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.

[35] Yaxing Wang, Salman Khan, Abel Gonzalez-Garcia, Joost van de Weijer, and Fahad Shahbaz Khan. Semi-supervised learning for few-shot image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4453–4462, 2020.

[36] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003.

[37] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE international conference on computer vision*, pages 1476–1485, 2019.

[38] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International Conference on Machine Learning*, pages 7354–7363. PMLR, 2019.

[39] Richard Zhang. Making convolutional networks shift-invariant again. In *ICML*, 2019.

[40] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

[41] Richard Zhang, Jun-Yan Zhu, Phillip Isola, Xinyang Geng, Angela S Lin, Tianhe Yu, and Alexei A Efros. Real-time

user-guided image colorization with learned deep priors. *ACM Transactions on Graphics (TOG)*, 9(4), 2017.

[42] Xuaner Zhang, Qifeng Chen, Ren Ng, and Vladlen Koltun. Zoom to learn, learn to zoom. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3762–3770, 2019.

[43] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 286–301, 2018.

[44] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.